

Ref #	Hits	Search Query	DBs	Default Operator	Plurals	Time Stamp
L1	510	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling relat\$3 lateral similar)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 19:33
L2	10	1 and (707/101 707/102 707/6). ccls.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:34
L3	7	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling sister brother lateral)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 19:34
L4	0	3 and (707/101 707/102 707/6). ccls.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:34
S27 9	15	(merg\$3 with (sibling relat\$3 lateral similar)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:34
S28 0	510	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling relat\$3 lateral similar)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:42
S28 1	7	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:38
S28 2	2413	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral)) and (usage)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:38
S28 3	83	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral)) and ((usage).clm.)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:38
S28 4	9273	usage and node and hierarch\$4	US-PGPUB; USPAT	OR	ON	2005/11/19 17:39
S28 5	2586	(usage and node and hierarch\$4) and (access\$3 with node)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:41
S28 6	50	((usage and node and hierarch\$4) and (access\$3 with node)) and (threshold with (usage access))) and merg\$3 and (split\$3 break\$3 divid\$3 division)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:40
S28 7	42	((usage and node and hierarch\$4) and (access\$3 with node)) and (threshold with (usage access))) and merg\$3 and (split\$3)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:40

S28 8	187	((usage and node and hierarch\$4) and (access\$3 with node)) and (threshold with (usage access))	US-PGPUB; USPAT	OR	ON	2005/11/19 17:40
S28 9	55	S280 and S285	US-PGPUB; USPAT	OR	ON	2005/11/19 17:43
S29 0	7	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral brother sister non-parent non-parent)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:43
S29 1	7	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral brother sister non-parent nonparent)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:43
S29 2	1	S291 and S285	US-PGPUB; USPAT	OR	ON	2005/11/19 17:44
S29 3	2	S291 and S284	US-PGPUB; USPAT	OR	ON	2005/11/19 17:44
S29 4	7	S291 and (usage)	US-PGPUB; USPAT	OR	ON	2005/11/19 17:44

Ref #	Hits	Search Query	DBs	Default Operator	Plurals	Time Stamp
L1	510	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling relat\$3 lateral similar)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 19:35
L2	10	1 and (707/101 707/102 707/6). ccls.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:34
L3	7	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling sister brother lateral)) and (usage with threshold)	US-PGPUB; USPAT	OR	ON	2005/11/19 19:36
L4	0	3 and (707/101 707/102 707/6). ccls.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:34
L7	1	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling sister brother lateral)) and (usage with threshold).clm.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:37
L8	1	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling lateral)) and (usage with threshold).clm.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:36
L9	85	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling sister brother lateral)) and (usage).clm.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:38
L10	2	((merg\$3 combin\$3 join\$3 consolidat\$3 integrat\$3 group\$3) with (sibling sister brother lateral)) and (usage) and threshold).clm.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:37
L12	0	9 and (707/6 707/101 707/102 706/20).ccls.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:39

Ref #	Hits	Search Query	DBs	Default Operator	Plurals	Time Stamp
L13	15	vishik.in.	US-PGPUB; USPAT	OR	ON	2005/11/19 19:39
L14	2	13 and usage	US-PGPUB; USPAT	OR	ON	2005/11/19 19:39



usage threshold merging sibling node

- 2002

Search

Ad  
Sch  
Sch

Scholar

Results 1 - 10 of about 67 for usage threshold merging sibling node. (0.01 seconds)

### Real-Time Data Access Control on B-Tree Index Structures

TW Kuo, CH Wei, KY Lam - PROC INT CONF DATA ENG, 1999 - [ieeexplore.ieee.org](#)... However, when splitting or **merging** of any **node** may be ... some bottleneck resource, eg, disks, ex- ceeds some **threshold**. ... the rest of this paper, our **usage** of the ...Cited by 9 - [Web Search](#) - [doi.ieeecomputersociety.org](#) - [doi.ieeeecs.org](#) - [csa.com](#) - [all 5 versions](#) »

### Adaptive Index Structures

Y Tao, D Papadias, CW Bay, H Kong - VLDB, 2002 - [vldb.org](#)... Since the pages allocated to **sibling** nodes are often not consecutive, a query (such as q ... to remedy this is to allocate several continuous pages to a **node** at a ...Cited by 6 - [View as HTML](#) - [Web Search](#) - [eden.dei.uc.pt](#) - [personal.psu.edu](#) - [repository.ust.hk](#) - [all 8 versions](#) »

### Multiversion Linear Quadtree for Spatio-Temporal Data

T Tzouramanis, M Vassilakopoulos, Y Manolopoulos - ADBIS-DASFAA, 2000 - [springerlink.com](#)... 7)) of the last inserted image: its **usage** is to ... If the latter number is below that **threshold**, then the ... of entries is above d), which is resolved by **merging** ...Cited by 13 - [Web Search](#) - [delab.csd.auth.gr](#) - [portal.acm.org](#) - [portal.acm.org](#)

### Real-Time Access Control and Reservation on B-Tree Indexed Data

TW Kuo, CH Wei, KY Lam - Real-Time Systems, 2000 - [kluweronline.com](#)... operations involving **node** splitting (and **merging**) are done ... of some bottleneck resource, eg, disks, exceeds some **threshold**. ... rest of this paper, our **usage** of the ...Cited by 1 - [Web Search](#) - [springerlink.com](#) - [cs.cityu.edu.hk](#) - [portal.acm.org](#) - [all 6 versions](#) »

### Automatic Web Page Classification in a Dynamic and Hierarchical Way

X Peng, B Choi - ICDM, 2002 - [doi.ieeecomputersociety.org](#)... is out of the range of **threshold** D, that ... After **merging** page vectors, the category information is changed ... profiles", KDD-99 Workshop on Web **Usage** Analysis and ...Cited by 6 - [Web Search](#) - [ieeexplore.ieee.org](#) - [benchoi.info](#) - [portal.acm.org](#) - [all 6 versions](#) »

### Subband HDTV Coding Using High-Order Conditional Statistics

SM Lei, TC Chen, KH Tzou - IEEE Journal on Selected Areas in Communications, 1993 - [ieeexplore.ieee.org](#)... Subsequently, code table **merging** can be done optimally ... coding rate (10) becomes less than a predetermined **threshold**. ... has a minor drawback in practi- cal **usage**. ...Cited by 1 - [Web Search](#) - [ieeexplore.ieee.org](#) - [csa.com](#)

### Merging Thesauri: Principles and Evaluation

H Mili, R Rada - IEEE Transactions on Pattern Analysis and Machine ..., 1988 - [doi.ieeeecs.org](#)... MILI AND RADA: **MERGING** THESAURI 207 ... stops whenever the second **node** is reached. ... by a Boolean query, once the size of the set is smaller than a given **threshold**. ...Cited by 11 - [Web Search](#) - [doi.ieeecomputersociety.org](#) - [portal.acm.org](#) - [ieeexplore.ieee.org](#) - [all 7 versions](#) »

### Visual Focusing and Transition Techniques in a Treeviewer for Web Information Access

K Wittenburg, E Sigman - VL, 1997 - [ieeexplore.ieee.org](#)... [20] have looked at **merging** bookmark files ... manipulations of the highlighting function based on community **usage**. ... hits up to a config- urable **threshold** appear in ...Cited by 11 - [Web Search](#) - [ieeexplore.ieee.org](#) - [portal.acm.org](#) - [portal.acm.org](#)

[PS] [An Investigation into Compositional Features and Feature Merging for Maximum Entropy-Based Parse ...](#)

T Mullen - 2002 - [iccs.informatics.ed.ac.uk](#)

... The discussion focuses on the difference between the cases in which **merging** might be of ... a telescope" is "at- tached" to the noun phrase (NP) **node** of the ...

[Cited by 1](#) - [View as HTML](#) - [Web Search](#) - [cogsci.ed.ac.uk](#) - [iccs.inf.ed.ac.uk](#) - [dspace.ub.rug.nl](#) - [all 5 versions](#) »

[A comparison of data \*\*merging\*\* methodologies for extending a microsimulation model](#)

D Schofield, J Polette, J McCrae, C O'Donoghue, M ... - 1996 - [natsem.canberra.edu.au](#)

... Figure 4 Stages in predicting child care **usage** ... Children in **sibling** or step care ... A

Comparison of Data **Merging** Methodologies for Extending a Microsimulation ...

[Cited by 1](#) - [View as HTML](#) - [Web Search](#) - [beje.decon.ufpe.br](#) - [papers.ssrn.com](#) - [ideas.repec.org](#) - [all 8 versions](#)

»

Google

Result Page:    1 2 3 4 5 6 7    **Next**

usage threshold merging sibling nod

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google



usage threshold merging sibling node

- 2002

Search

Ad  
Sch  
Sch

Scholar

Results 11 - 20 of about 67 for usage threshold merging sibling node. (0.01 seconds)

[PS] Designing access methods for bitemporal databases

A Kumar, VJ Tsotras, C Faloutsos - IEEE Transactions on Knowledge and Data Engineering, 1998 - cs.ucr.edu  
... the Bitemporal R-Tree based on various **merging** and splitting ... The **usage** of two R-trees  
reminds the Dual ... **threshold**, a vacuuming process completely vacuums all its ...

Cited by 104 - [View as HTML](#) - [Web Search](#) - [portal.acm.org](#) - [ieeexplore.ieee.org](#) - [drum.umd.edu](#) - [all 9 versions](#)

»

Quadtree-structured variable-size block-matching motion estimation with minimal error

I Rhee, GR Martin, S Muthukrishnan, RA Packwood - IEEE TRANS CIRCUITS SYST VIDEO TECHNOL, 2000 -  
ieeexplore.ieee.org

... the number of blocks if the bit **usage** of each ... the resulting error is above a prescribed  
**threshold**, then that ... Finally, a process of re-**merging** small blocks to ...

Cited by 7 - [Web Search](#) - [ieeexplore.ieee.org](#) - [csa.com](#)

Variable n-grams and extensions for conversational speech language modeling

M Siu, M Ostendorf - IEEE Transactions on Speech and Audio Processing, 2000 - ieeexplore.ieee.org

... model represents prior knowledge of word **usage** in a ... a **threshold** on n-gram occurrence  
counts to build ... A. **Node Merging** Previous variable n-gram work combines a ...

Cited by 21 - [Web Search](#) - [ieeexplore.ieee.org](#)

A Comparison of Access Methods for Temporal Data

B Salzberg, VJ Tsotras - ACM Computing Surveys, 1999 - cs.auc.dk

Page 1. -1 A Comparison of Access Methods for Temporal Data Betty Salzberg and Vassilis

J. Tsotras June 13, 1997 TR-18 A TimeCenter Technical Report Page 2. 0 ...

Cited by 39 - [View as HTML](#) - [Web Search](#)

Parallel volume rendering using binary-swap compositing

KL Ma, JS Painter, CD Hansen, MF Krogh - IEEE Computer Graphics and Applications, 1994 -  
ieeexplore.ieee.org

... for volume rendering if the memory **usage** on each ... or until the accumulated opacity  
reaches a **threshold** cut-off ... is a naive approach for parallel **merging** of the ...

Cited by 101 - [Web Search](#) - [portal.acm.org](#) - [portal.acm.org](#) - [csa.com](#) - [all 6 versions](#) »

Lexical acquisition with WordNet and the Mikrokosmos ontology

T O'Hara, K Mahesh, S Nirenburg - COLING/ACL Wordkshop on **Usage** of WordNet in NLP Systems, 1998 -  
acl.ldc.upenn.edu

... a dictionary, it provides definitions and **usage** examples ... Plus, since the children  
& **sibling** matches are ... match, provided the score is above a certain **threshold**. ...

Cited by 8 - [View as HTML](#) - [Web Search](#) - [acl.eldoc.ub.rug.nl](#) - [ai.sri.com](#) - [cs.nmsu.edu](#) - [all 5 versions](#) »

[PS] Parallel volume rendering using binary-swap image composition

KL Ma, JS Painter, CD Hansen, MF Krogh - IEEE Computer Graphics and Applications, 1994 - cis.ohio-state.edu

... volume rendering provided that the memory **usage** on each ... or until the accumulated  
opacity reaches a **threshold** cut-o ... A naive approach for parallel **merging** of the ...

Cited by 18 - [View as HTML](#) - [Web Search](#) - [cs.ucdavis.edu](#) - [cse.ohio-state.edu](#)

A Framework for Dependability Driven Software Integration

N Suri, S Ghosh, T Marlowe - The 1998 18 th International Conference on Distributed ..., 1998 - doi.ieeeecs.org

... will use the terms parent, child, and **sibling** in the ... R3: Future integration by **merging**:

An FCM can be integrated ... can be measured from previ- ous **usage** of that ...

Cited by 2 - [Web Search](#) - [ieeexplore.ieee.org](#) - [cs.virginia.edu](#) - [deeds.informatik.tu-darmstadt.de](#) - [all 12 versions](#) »

### Enhanced Automatic Creation of Multi-Purpose Object Hierarchies

J Haber, M Stamminger, HP Seidel - Pacific Conference on Computer Graphics and Applications, 2000 - doi.ieeecomputersociety.org

... The diameter of the root **node** is 8, so the cost of the left, flat hierarchy ... sort the children Ó with respect to the costs that would arise by **merging** Ó with ...

Cited by 2 - [Web Search](#) - [doi.ieee.org](#) - [ieeexplore.ieee.org](#) - [portal.acm.org](#) - [all 6 versions](#) »

### Survey of clustering data mining techniques

P Berkhin - Accrue Software, 2002 - cs.dal.ca

... cluster **node** contains child clusters; **sibling** clusters partition ... smaller than a certain **threshold**, by using ... and proceeds iteratively with **merging** or splitting ...

Cited by 116 - [View as HTML](#) - [Web Search](#) - [it.bond.edu.au](#) - [cs.itu.edu.tr](#) - [datanautics.com](#) - [all 23 versions](#) »



Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [Next](#)

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google





usage threshold merging sibling node

- 2002

Search

Ad  
Sci  
Sci

Scholar

Results 21 - 30 of about 67 for usage threshold merging sibling node. (0.02 seconds)

### LDI Tree: A Hierarchical Representation for Image-Based Rendering

CF Chang, G Bishop, A Lastra - SIGGRAPH, 1999 - portal.acm.org

... output LDI, which is similar to the pixel **merging** in the ... The difference is that a single **threshold** value of depth ... The memory **usage** of the LDI trees is shown in ...Cited by 73 - [Web Search](#) - [cse.psu.edu](#) - [cs.princeton.edu](#) - [cs.unc.edu](#) - [all 12 versions](#) »

### Reconstruction of vascular networks using three-dimensional models

P Hall, M Ngan, P Andrae - IEEE Transactions on Medical Imaging, 1997 - ieeexplore.ieee.org

... it to be fast, reasonable in memory **usage**, and reliable ... vasculature is added to the VC by **merging** it with ... If this falls below some **threshold**, the sliver-path ...Cited by 14 - [Web Search](#) - [ieeexplore.ieee.org](#) - [ncbi.nlm.nih.gov](#)

### Algorithmic Geometry

JD Boissonnat, JD Boissonnat, JD Boissonnat - 1998 - print.google.com

... We have followed the French text in systematically using the words **node** and arc ... text for the word saillant (meaning salient) to follow the **usage** with convex ...Cited by 126 - [Web Search](#) - [inria.fr](#) - [all 3 versions](#) » - [Library Search](#)

### Discovering user communities on the Internet using unsupervised machine learning techniques

G Paliouras, C Papatheodorou, V Karkaletsis, CD ... - Interacting with Computers, 2002 - delos.di.uoa.gr

... two clusters into a new one (**merging**) and dividing ... a full investigation of the effect of the connectivity **threshold**. ... for only a small part of the **usage** of the ...Cited by 8 - [View as HTML](#) - [Web Search](#) - [iit.demokritos.gr](#) - [ingentaconnect.com](#) - [all 6 versions](#) »

### A Survey of Indexing Techniques for Semistructured Documents

F Weigel, PR des Fortgeschrittenenpraktikums - Institute for Computer Science, University of Munich, 2002 - cis.uni-muenchen.de

... 33 5.5 **Node** identification . . . . . 35 5.5.4 **Sibling** numbers . . . . .Cited by 4 - [View as HTML](#) - [Web Search](#) - [postech.ac.kr](#) - [pms.informatik.uni-muenchen.de](#) - [pms.ifi.lmu.de](#) - [all 7 versions](#) »

### Comparison of access methods for time-evolving data

B Salzberg, VJ Tsotras - ACM Computing Surveys, 1999 - portal.acm.org

... for (a) locating the relative position of a timestamp on the evolution of a given branch and (b) locating the same timestamp among **sibling** branches. ...Cited by 143 - [Web Search](#) - [ad.informatik.uni-freiburg.de](#) - [ccs.neu.edu](#) - [sok.susu.ru](#) - [all 5 versions](#) »

### Heap implementations and variations

J Bojesen - Written Project, Department of Computing, University of ..., 1998 - dimcom.uqac.quebec.ca

... 29 7.2 **Node** rell strategy . . . . . 29 ...Cited by 3 - [View as HTML](#) - [Web Search](#) - [diku.dk](#)

### Searching the Web

A Arasu, J Cho, H Garcia-Molina, A Paepcke, S ... - ACM Transactions on Internet Technology, 2001 - portal.acm.org

... Crawl control may also use feedback from **usage** patterns to guide the crawling process(connection between the query ... Crawl & Stop with **Threshold**: We again ...Cited by 130 - [Web Search](#) - [snoopy.cs.nccu.edu.tw](#) - [utdallas.edu](#) - [dsi.uniroma1.it](#) - [all 38 versions](#) »

[PS] [Towards multi-domain speech understanding with flexible and dynamic vocabulary](#)

GYC Chung - 2001 - [sls.lcs.mit.edu](#)

Page 1. Towards Multi-Domain Speech Understanding with Flexible and Dynamic

Vocabulary by Grace Chung SM, Massachusetts Institute ...

[Cited by 3](#) - [View as HTML](#) - [Web Search](#) - [language.cnri.reston.va.us](#) - [sls.csail.mit.edu](#) - [dspace.mit.edu](#) - [all 7 versions](#) »

[Visualization of Large Terrains Made Easy](#)

P Lindstrom, V Pascucci - IEEE Visualization, 2001 - [portal.acm.org](#)

... accuracy, mesh com- plexity, memory usage, refinement speed ... the projected errors

meet some threshold or the ... resolution mesh and recursively merging pairs of ...

[Cited by 67](#) - [Web Search](#) - [citeseer.csail.mit.edu](#) - [cc.gatech.edu](#) - [cgvr.korea.ac.kr](#) - [all 20 versions](#) »

◀ Google ▶

Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [Next](#)

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google



usage threshold merging sibling node

- 2002

Search

[Ad](#)  
[Sch](#)  
[Sch](#)
**Scholar**

Results 31 - 40 of about 67 for usage threshold merging sibling node. (0.02 seconds)

Taking a walk in a planar arrangement

S Har-Peled - ANNU SYMP FOUND COMPUT SCI PROC. pp. 100-110. 1999, 1999 - [epubs.siam.org](http://epubs.siam.org)  
 ... v). To perform in step (ii) the **merging** of the ... (namely, until we reach a **node** that  
 is ... v), the procedure that computes the "**sibling**" transient trapezoids ...

Cited by 25 - [Web Search](#) - [ieeexplore.ieee.org](http://ieeexplore.ieee.org) - [doi.ieeecs.org](http://doi.ieeecs.org) - [uiuc.edu](http://uiuc.edu) - [all 17 versions](#) »

Database systems for structured documents

R Sacks-Davis, T Arnold-Moore, J Zobel - IEICE Transactions on Information and Systems, 1995 -  
[mds.rmit.edu.au](http://mds.rmit.edu.au)

Page 1. Database Systems for Structured Documents Ron Sacks-Davis Timothy  
 Arnold-Moore y Justin Zobel z Abstract Documents stored ...

Cited by 51 - [View as HTML](#) - [Web Search](#) - [pms.informatik.uni-muenchen.de](http://pms.informatik.uni-muenchen.de) - [cs.rmit.edu.au](http://cs.rmit.edu.au) - [csa.com](http://csa.com)

Tailor: a layout system based on trapezoidal corner stitching

D Marple, M Smulders, H Hegen, D Vangheluwe - IEEE Transactions on Computer-Aided Design of Integrated  
 ..., 1990 - [ieeexplore.ieee.org](http://ieeexplore.ieee.org)

... The diffusion contacts appear in both layout planes. N = blat uuirrber of objects  
 o = objects in area/shadow/node 'l = **threshold** uumbec MARPLE crc!! ...

Cited by 22 - [Web Search](#) - [ieeexplore.ieee.org](http://ieeexplore.ieee.org) - [csa.com](http://csa.com)

[PS] Basic External Memory Data Structures

R Pagh - Algorithms for Memory Hierarchies, 2002 - [it-c.dk](http://it-c.dk)

... larger than ordinary pointers) the space used for implementing robust pointers  
 increases total space **usage** by at ... **node** is the distance to its descendant leaves. ...

Cited by 1 - [View as HTML](#) - [Web Search](#)

[book] Real World Semantic Web Applications

V Kashyap, V Kashyap, L Shklar - 2002 - [print.google.com](http://print.google.com)

... above standards, though new advancements may see re-alignments and **merging** of the ...  
 ecosystem" for the Semantic Web will form, lowering the **threshold** for entry ...

[Web Search](#) - [Library Search](#)

Out-Of-Core Algorithms for Scientific Visualization and Computer Graphics

CT Silva, YJ Chiang, J El-Sana, P Lindstrom - Visualization'02, 2002 - [cs.utah.edu](http://cs.utah.edu)

... memory for each sub-list that is larger than main memory in the above **merging** step ...  
 Each tree **node** corresponds to one disk block, capable of holding up to B items ...

Cited by 13 - [View as HTML](#) - [Web Search](#) - [cse.ogi.edu](http://cse.ogi.edu) - [citeseer.csail.mit.edu](http://citeseer.csail.mit.edu) - [cc.gatech.edu](http://cc.gatech.edu) - [all 7 versions](#) »

[PS] Modular Stochastic HPSGs for Question Answering

V Keselj - 2002 - [cs.dal.ca](http://cs.dal.ca)

... memory management within the algorithm, sub-**node** hid- ... Modularity in grammar **usage**  
 includes the problems of retrieving and **merging** several grammars into one at ...

Cited by 2 - [View as HTML](#) - [Web Search](#)

A Nearly Optimal Deterministic Parallel Voronoi Diagram Algorithm

R Cole, MT Goodrich, CO Dunlaing - Algorithmica, 1996 - [springerlink.com](http://springerlink.com)

... (c) **Usage** of duplicate ... parallel operations on arrays, including Valiant's **merging**  
 technique and ... without cycles; it does not have a distinguished root **node**. ...

Cited by 3 - [Web Search](#)

Application of the SEPA Methodology and Tool Suite to the National Cancer Institute

KS Barber, TJ Graser, SR Jernigan, BJ McGiverin, J ... - HICSS, 1999 - dx.doi.org

... While recognizing the need for a unified domain model, DSSA does not outline a formal method for merging the requirements from multiple domain experts. ...

Cited by 10 - [Web Search](#) - [ieeexplore.ieee.org](#) - [doi.ieeecomputersociety.org](#) - [computer.org](#) - [all 12 versions »](#)bookj Advanced Database Indexing

Y Manolopoulos, Y Theodoridis, VJ Tsotras - 2000 - print.google.com

... 3 1 Figure 2.5. Combining merging steps. 3 1 ... 132 Figure 6. 1 1 . Object o's MBR meets q's MBR while the covering node rectangle N does not. 135 ...

Cited by 44 - [Web Search](#) - [kap.nl](#) - [cs.ucr.edu](#) - [all 6 versions »](#)Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [Next](#) [Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google



usage threshold merging sibling node

- 2002

Search

Adv  
Sch  
Sch**Scholar**

Results 41 - 50 of about 67 for usage threshold merging sibling node. (0.02 seconds)

bookj Geographical Information: From Research to Application Through Cooperation

M Rumor, R McMillan, HFL Ottens - 1996 - print.google.com

Page 1. GEOGRAPHICAL INFORMATION FROM RESEARCH TO APPLICATION THROUGH COOPERATION

Th isOne I II WJW8-1DZ-A974 Page 2. Geographical Information ...

Cited by 1 - Web SearchCross-Layer Design for Power Aware Communication in Multihop Wireless Networks

Y Xue - Master's Thesis, Department of Computer Science, University ..., 2002 - princeton.edu

... A The received power at # must at least equal to the minimum received power **threshold**,£¥¤ % & j ... is the total noise that **node** # observes on the ...View as HTML - Web Search - cairo.cs.uiuc.eduMillipede: Easy Parallel Programming in Available Distributed Environments

A ITZKOVITZ, A SCHUSTER - SOFTWARE—PRACTICE AND EXPERIENCE, 1997 - doi.wiley.com

... sorting the two parts concurrently (using pparblock), and sequentially **merging** thevalues ... One of each pair of **sibling** activities will always update the variable ...Web Search - doi.wiley.comARMS: an algebraic recursive multilevel solver for general sparse linear systems

Y Saad, B Suchomel - Numerical Linear Algebra with Applications, 2002 - doi.wiley.com

... preconditioned GMRES [13] as well as its scalar **sibling**, ILUM. For certain hardproblems, these attributes come with the added benefit of smaller memory **usage**. ...Cited by 36 - Web Search - cs.umn.edu - cs.umn.edu - www-users.cs.umn.edubookj Web prefetching using partial match prediction

T Palpanas, A Mendelzon - 2000 - cs.toronto.edu

... derived from its **usage** [PM94, NGBS + 97]. ... les, while a directed weighted arcfrom **node** A to **node** B expresses a metric for the ... **sibling** page. ...Cited by 47 - View as HTML - Web Search - cs.ucr.edu - workshop99.ircache.net - all 9 versions » - Library SearchA New Virtual Memory Implementation for L4/MIPS

C Szmajda - University of New South Wales, 1999 - itk.ntnu.no

... L4/MIPS data structures are also re-designed for efficiency and reduced memory **usage**. ...**Node** Joining and Doubling ... 6.4 **Merging** Flush Lists into Page Table Entries ...Cited by 3 - View as HTML - Web Search - disy.cse.unsw.edu.au - cse.unsw.edu.auAnatomy of a native XML base management system

T Fiebig, S Helmer, CC Kanne, G Moerkotte, J ... - VLDB Journal, 2002 - portal.acm.org

... guarantees serializability even if transactions directly access some **node** in a ... XML,which are related to the different **usage** profiles (coarse ... as "**threshold**". ...Cited by 45 - Web Search - springerlink.com - csd.uch.gr - pi3.informatik.uni-mannheim.de - all 9 versions »[P5] The Deevolution of Concurrent Logic Programming Languages

E Tick - JLP, 1995 - cir1.uoregon.edu

Page 1. The Deevolution of Concurrent Logic Programming Languages Evan

Tick University of Oregon CIS-TR{94{07 March 1994 Abstract ...

Cited by 25 - View as HTML - Web Search

SilkRoute: A framework for publishing relational data in XML

M Fernandez, Y Kadiyska, D Suciu, A Morishima, WC ... - ACM Transactions on Database Systems, 2002 - portal.acm.org

... The finest decomposition (each **node** is a separate partition) results in multiple ...

The only data manipulation performed by SilkRoute is **merging** of sorted tuple ...

Cited by 60 - [Web Search](#) - [cs.washington.edu](#) - [pages.stern.nyu.edu](#) - [soe.ucsc.edu](#) - [all 10 versions »](#)

Robust automated topic identification

CY Lin - 1997 - isi.edu

... R t . . . . . 35 3.4.2.1 Selecting Interesting Concepts Using Branch Ratio

**Threshold** . . . . . 59 3.5.5 Counting and Merging Techniques . . .

Cited by 19 - [View as HTML](#) - [Web Search](#) - [isi.edu](#) - [portal.acm.org](#) - [portal.acm.org](#)



Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [Next](#)

usage threshold merging sibling nod

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google



usage threshold merging sibling node

- 2002

Search

[Adv](#)  
[Sch](#)  
[Sch](#)
**Scholar**

Results 51 - 60 of about 67 for usage threshold merging sibling node. (0.02 seconds)

Automatic Construction of News Hypertext

T Dalamagas, MD Dunlop - HIM, 1997 - dbnet.ece.ntua.gr

... This suggests the **usage** of an inverse document frequency ... Documents that their similarity value with the query is above a predefined **threshold** are considered to ...Cited by 5 - [View as HTML](#) - [Web Search](#) - [ls1-www.cs.uni-dortmund.de](#) - [ls1-www.informatik.uni-dortmund.de](#) - [cs.strath.ac.uk](#) - [all 9 versions](#) »Subword lexical modelling for speech recognition

R Lau, S Seneff, E Grimson, V Zue, M Cambridge - 1998 - groups.csail.mit.edu

... 94 5.3.3 **Merging** Theories to Increase Pruning Opportunities . . .

. 95 5.3.4 Syllables vs. Words . . . . .

Cited by 12 - [View as HTML](#) - [Web Search](#) - [raylau.com](#) - [sls.csail.mit.edu](#) - [portal.acm.org](#) - [all 9 versions](#) »LDI tree: a sampling rate preserving and hierarchical data representation for image-based rendering

CF Chang, C Hill, N England, RMT II - 2001 - cs.unc.edu

... 2 Sample **Merging** and Redundancy ... 40 4.1 Memory **Usage** .....40 ...[View as HTML](#) - [Web Search](#) - [cs.unc.edu](#)[PS] FeasPar- A Feature Structure Parser Learning to Parse Spontaneous Speech

FD Bu - 1996 - archive.cis.ohio-state.edu

... Since the entire state of a parse only involves three pieces of information (current position, current **node**, and return points), backtracking and search is ...Cited by 9 - [View as HTML](#) - [Web Search](#) - [funet.fi](#) - [funet.fi](#) - [is.cs.cmu.edu](#) - [all 5 versions](#) »[PS] LR parsing for Tree Adjoining Grammars and its application to corpus-based natural language parsing

CA Prolo - 2002 - cis.upenn.edu

... 135 3.37 Possible positions for unmarked nodes relative to the dotted **node** . . . . .186 5.11 **Usage** of the QP label in the PTB . . . . . 195 5.20 **Merging** labels . . . . .Cited by 4 - [View as HTML](#) - [Web Search](#) - [repository.upenn.edu](#) - [repository.upenn.edu](#)[PS] Lexical Chains for Summarization

R Barzilay - 1997 - sls.lcs.mit.edu

... We present a new algorithm to compute lexical chains in a text, **merging** several robust knowledge sources ... The **usage** of "this measure" to the "error analysis" ...Cited by 12 - [View as HTML](#) - [Web Search](#) - [sls.csail.mit.edu](#) - [cs.cornell.edu](#)[PS] Incremental algorithms for general purpose layout system

MH Cynn - 1994 - historical.ncstrl.org

... using non-recursive area enumeration. The system consists of three stages:

**node** extraction, ... 56 4 Plane Generation and **Node** Extraction ...[View as HTML](#) - [Web Search](#) - [cs.uiuc.edu](#) - [Library Search](#)[BOOK] Resource and knowledge discovery from the Internet and multimedia repositories

R Zaiane - 2001 - fas.sfu.ca

... 20 2.1.3 Web **Usage** Mining . . . . . (See Figures A.6 and A.7 in Appendix A). Users may have quite different backgrounds, interests, and purposes of **usage**. ...

[Cited by 12](#) - [View as HTML](#) - [Web Search](#) - [fas.sfu.ca](#) - [cs.sfu.ca](#) - [all 6 versions »](#) - [Library Search](#)

[bookj](#) [Metadata enhanced content management in media companies](#)

S Jokela - 2001 - [lib.tkk.fi](#)

... of the main factors affecting production, delivery, and **usage** of content is convergence:

Communications, computing, and content industries are **merging** into a ...

[Cited by 6](#) - [View as HTML](#) - [Web Search](#) - [lib.hut.fi](#) - [Library Search](#)

[IMPS: implicit surfaces for interactive animated characters](#)

KB Russell - 1999 - [xenia.media.mit.edu](#)

... 69 A.2 **Usage** Notes and Examples ... 78 6 Page 7. B.3

Developing New **Node** Types : : : : 79 ...

[View as HTML](#) - [Web Search](#) - [vismod.www.media.mit.edu](#) - [www-white.media.mit.edu](#) - [media.mit.edu](#) - [all 7 versions »](#)



Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#) [Next](#)

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google





usage threshold merging sibling node

- 2002

Search

[Adv](#)  
[Sch](#)  
[Sch](#)
**Scholar**

Results 61 - 67 of 67 for usage threshold merging sibling node. (0.03 seconds)

Portable, modular expression of locality

DP Stoutamire, JA Feldman - 1997 - icsi.berkeley.edu

... The write buffer does write-**merging** if it is capable of combin- ing multiple writes to the same memory region without requiring additional transactions with ...

Cited by 8 - [View as HTML](#) - [Web Search](#) - [icsi.berkeley.edu](#) - [david.stoutamire.com](#) - [portal.acm.org](#) - [all 9 versions](#) »

[book] Mobile Commerce: Opportunities, Applications, and Technologies of Wireless Business

P May - 2001 - print.google.com

... of their teachers and parents, and habituating themselves for a lifetime of mobile service **usage**. ... with each other, the addition of a single new **node** has a high ...

Cited by 24 - [Web Search](#) - [Library Search](#)

A cost-benefit approach to resource allocation in scalable metacomputers

RS Borgstrom, B Awerbuch, Y Amir - 2001 - cnds.jhu.edu

... This approach gives us an online strategy provably competitive with the optimal offline algorithm in the maximum **usage** of each resource. ...

Cited by 3 - [View as HTML](#) - [Web Search](#) - [cnds.jhu.edu](#) - [csa.com](#) - [Library Search](#)

The art of computer programming, volume 3:(2nd ed.) sorting and searching

DE Knuth - 1998 - portal.acm.org

... Zhu Hong , Robert Sedgewick, Notes on **merging** networks (Preliminary Version ... J. Piestrak, The Minimal Test Set for Multioutput **Threshold** Circuits Implemented as ...

Cited by 16 - [Web Search](#) - [portal.acm.org](#)

Linux Advanced Routing & Traffic Control HOWTO

B Hubert, G Maxwell, R van Mook, M van Oosterhout, ... - setembro de, 2002 - ds9a.nl

Page 1. Linux Advanced Routing &amp; Traffic Control HOWTO Bert Hubert Netherlabs BV

bert.hubert@netherlabs.nl Thomas Graf (Section Author) tgraf@suug.ch ...

Cited by 36 - [View as HTML](#) - [Web Search](#) - [dug.net.pl](#) - [lartc.org](#) - [sti.uniurb.it](#) - [all 93 versions](#) »

Practical parallel divide-and-conquer algorithms

JC Hardwick, A Beguelin, B Maggs - 1997 - www-2.cs.cmu.edu

... 90 6.1.1 Choice of load-balancing **threshold** . . . . . 93 6.3 Effect of

load-balancing **threshold** on quicksort on the Cray T3D . . . . .

Cited by 3 - [View as HTML](#) - [Web Search](#) - [citeseer.csail.mit.edu](#) - [csa.com](#) - [all 7 versions](#) » - [Library Search](#)

Design dependencies within the automatic generation of hypermedia presentations

OR Martinez - Master's thesis, Technical University of Catalonia, June, 2002 - cwi.nl

... Presen- tation System (IMMPS) [43], is a difficult task that requires **merging** techniques from ... key of multimedia is not in a simple and plain **usage** of different ...

Cited by 5 - [View as HTML](#) - [Web Search](#) - [homepages.cwi.nl](#) - [Library Search](#)

Result Page: [Previous](#) [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#)

usage threshold merging sibling nod

[Google Home](#) - [About Google](#) - [About Google Scholar](#)

©2005 Google

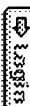


IEEE COMPUTER SOCIETY

Search:

Go

Advanced Search

[Home](#)
[Digital Library](#)
[Site Map](#)
[Store](#)
[Help](#)
[Contact Us](#)
[Press Room](#)
[Shopping Cart](#)


## digital library

## DIGITAL LIBRARY HOME

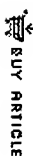
[BROWSE BY TITLE](#)[BROWSE BY SUBJECT](#)[SEARCH](#)[LIBRARY/INSTITUTION  
RESOURCES](#)[RESOURCES](#)[SUBSCRIPTION](#)[ABOUT THE DIGITAL LIBRARY](#)[Past Issues >> Table of Contents >> Abstract](#)
**IEEE TRANSACTIONS ON  
KNOWLEDGE AND  
DATA ENGINEERING**

May 2005 (Vol. 17, No. 5) pp. 614-627

**WISDOM: Web Intrapage Informative Structure  
Mining Based on Document Object Model**

Hung-Yu Kao  
Jan-Ming Ho, IEEE  
Ming-Syan Chen, IEEE

Full Article Text:

DOI Bookmark: <http://doi.ieeecomputersociety.org/10.1109/TKDE.2005.84>**Abstract**

To increase the commercial value and accessibility of pages, most content sites tend to publish their pages with intrasite redundant information, such as navigation panels, advertisements, and copyright announcements. Such redundant information increases the index size of general search engines and causes page topics to drift. In this paper, we study the problem of mining intrapage informative structure in news Web sites in order to find and

## Abstract Contents:

Abstract  
Index Terms  
Citation

## Free access to

- ☐ Abstracts
- ☐ Selected PDFs

Electronic subscribers log  
in to

- ☐ Access HTML/PDFs of full text articles
- ☐ Download full issue (ZIP of PDFs)

## Subscription information

Get a Web account

eliminate redundant information. Note that intrapage informative structure is a subset of the original Web page and is composed of a set of fine-grained and informative blocks. The intrapage informative structures of pages in a news Web site contain only anchors linking to news pages or bodies of news articles. We propose an intrapage informative structure mining system called WISDOM (Web Intrapage Informative Structure Mining based on the Document Object Model) which applies Information Theory to DOM tree knowledge in order to build the structure. WISDOM splits a DOM tree into many small subtrees and applies a top-down informative block searching algorithm to select a set of candidate informative blocks. The structure is built by expanding the set using proposed merging methods. Experiments on several real news Web sites show high precision and recall rates which validates WISDOM's practical applicability.

---

#### Additional Information

[Back to Top](#)

---

**Index Terms**- Intrapage informative structure, DOM, entropy, information extraction.

**Citation:** Hung-Yu Kao, Jan-Ming Ho, Ming-Syan Chen. "WISDOM: Web Intrapage Informative Structure Mining Based on Document Object Model," *IEEE Transactions on Knowledge and Data Engineering*, vol. 17, no. 5, pp. 614-627, May, 2005.

---

Usage of this product signifies your acceptance of the Terms of Use.

This site and all contents (unless otherwise noted) are Copyright © 2005, IEEE, Inc. All rights reserved.

# WISDOM: Web Intrapage Informative Structure Mining Based on Document Object Model

Hung-Yu Kao, Jan-Ming Ho, *Member, IEEE*, and Ming-Syan Chen, *Fellow, IEEE*

**Abstract**—To increase the commercial value and accessibility of pages, most content sites tend to publish their pages with intrasite redundant information, such as navigation panels, advertisements, and copyright announcements. Such redundant information increases the index size of general search engines and causes page topics to drift. In this paper, we study the problem of mining *intrapage informative structure* in news Web sites in order to find and eliminate redundant information. Note that intrapage informative structure is a subset of the original Web page and is composed of a set of fine-grained and informative blocks. The intrapage informative structures of pages in a news Web site contain only anchors linking to news pages or bodies of news articles. We propose an intrapage informative structure mining system called *WISDOM (Web Intrapage Informative Structure Mining based on the Document Object Model)* which applies Information Theory to DOM tree knowledge in order to build the structure. WISDOM splits a DOM tree into many small subtrees and applies a top-down informative block searching algorithm to select a set of candidate informative blocks. The structure is built by expanding the set using proposed merging methods. Experiments on several real news Web sites show high precision and recall rates which validates WISDOM's practical applicability.

**Index Terms**—Intrapage informative structure, DOM, entropy, information extraction.

## 1 INTRODUCTION

MANY Web pages are generated online for Web site maintenance, flexibility, and scalability purposes. They are usually generated by putting page content stored in back-end databases into predefined templates. The experimental results in [4] show that, on average, 43 percent of Web pages contain templates which indicates how pervasive template usage has become. Most commercial Web sites, such as search engines, portal sites, e-commerce stores, and news, apply a systematic technique to generate Web pages and to adapt various requests from numerous Web users. These sites are referred to as *systematic Web sites* [16]. The evolution of automatic Web page generation and the sharp increase of systematic Web sites have contributed to the explosive growth of Web page numbers. There exists much redundant and irrelevant information in these Web pages [1], [23], such as navigation panels, advertisements, catalogs of services, and announcements of copyright and privacy policies which are distributed over almost all pages of a systematic Web site. Such information is still crawled and indexed by search engines and information agents, thus significantly increasing corresponding storage and computing overhead.

We define specific regions of a page that users are interested in as *informative blocks* (or referred to as *IB*). Information within IBs manifests the main topic of the page

and indicates related information. The set of these blocks and corresponding connecting structures form the *informative structure* (or referred to as *IS*) of the page. Fig. 1 shows the *IS* of an example news page and its corresponding parts of content. A Web page can be represented by a tree structure, i.e., Document Object Model (DOM) [27] and each content block in a page is a subtree of the original DOM tree. The *IS* can be defined as a reduced tree united by subtrees of *IBs*. The tree relation of united subtrees is also kept in *IS*. The *IS*, for example, in Fig. 1 is built by uniting subtrees of two news table of content blocks and is a tree reduced from the original DOM tree. As proposed in [16], which deals with the *IS* in a site, called *interpage informative structure*. The structure is composed of informative pages within a Web site and interconnecting links. In this paper, we work with *ISs* of individual pages, called *intrapage informative structure* (for simplicity, we use the same denotation *IS* in the last parts of the paper). Each page has its own *IS* and the structure is composed of *IBs* within the page.

Web informative content mining is an important task for search engines and Web agents [11]. Internet crawlers can use the *IS* to focus on crawling informative paths. Search engines can reduce the size of indices and make them more precise by removing the redundant and irrelevant page blocks. Intermedia information agents that search for specific information among Web sites with different presentation styles, page layouts, and site mapping can also benefit from the information preprocessed by the structure.

News search engines like Google News, Altavista News, and NSE<sup>1</sup> are typical examples of intermedia information agents. They crawl diverse news articles from thousands of news Web sites and extract and index article blocks. The *IS* of a page consists of sets of table of contents (abbreviated as *TOC*) blocks and news article blocks. The structure helps

- H.-Y. Kao is with the Department of Computer Science and Information Engineering, National Cheng Kung University, Tainan, Taiwan, ROC. E-mail: hykao@mail.ncku.edu.tw.
- J.-M. Ho is with the Institute of Information Science, Academia Sinica, Taipei, Taiwan, ROC. E-mail: hoho@iis.sinica.edu.tw.
- M.-S. Chen is with the Graduate Institute of Communication Engineering and the Electrical Engineering Department, National Taiwan University, Taipei, Taiwan, ROC. E-mail: mschen@cc.ee.ntu.edu.tw.

Manuscript received 8 Sept. 2003; revised 3 June 2004; accepted 22 Sept. 2004; published online 17 Mar. 2005.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number TKDE-0180-0903.

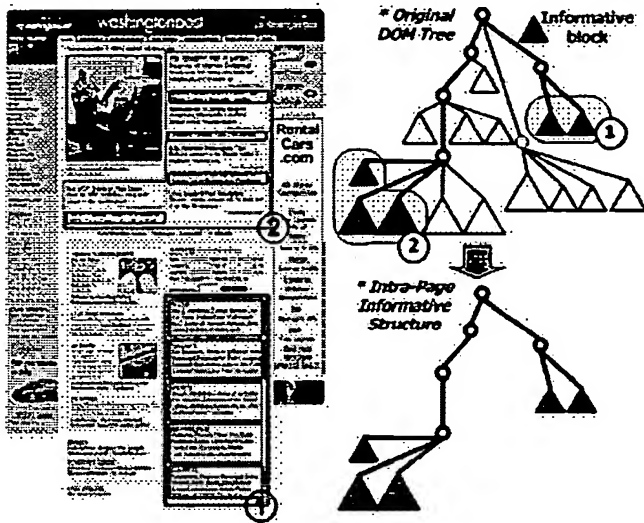


Fig. 1. The original document structure and its IS of a news page from WashingtonPost.

agents to automatize crawling and indexing. The IS of article pages usually consists of fine-grained and joined IBs and most of them contain only one HTML tag in its subtree, i.e., leaf nodes in the DOM tree. These blocks contain news article metadata, such as title, date, reporter, and place, which are very useful in categorizing pages for news information agents and metadata extraction. Agents can automatically extract article metadata by using this structure.

According to the definitions of hubs and authorities in [18], a good hub is a page linking to a good authority page that is relevant to some specific query. Analogously, we define a good *information hub* as a block linking to a good *information authority* block which will provide useful information. The IS of a page can then be considered as the set of blocks of good information hubs and good information authorities within that page. Note that Web pages in news Web sites usually contain the obvious and clear ISs, i.e., TOC and article blocks, in our observation.

Fig. 2 shows the root page of the news Web site WashingtonPost (<http://www.washingtonpost.com>) and blocks 1 and 2 which provide anchors linking to hot news and selected news are the information hubs. We consider these two blocks as IBs as they are the crawling points for news information agents to collect daily news. Block 3 is merely a menu block which is appended ubiquitously to most pages in the WashingtonPost Web site, and is thus considered as redundant.

In an HTML document, tags are inserted for purposes of the page layout, content presentation, and for providing interactive functions, e.g., form filling and document linking. After being rendered by the browser, tags are invisible to users and are represented by means of visual appearances and functions. The layout and style of presentations provide hints to users for accessing and understanding information easily. The corresponding tagging structure therefore contains information about representation and semantics of Web pages. For example, a group of tight sibling anchor nodes with the short anchor-text, e.g., the tagging tree of block 3 in Fig. 2, is different from a group of sibling nodes with the long anchor-text interleaved with context nodes, e.g., the tagging tree of block 2 in Fig. 2. In Fig. 2, block 1 containing several tightly coupled anchor groups also provides different functionality and representation from block 2 and block 3. In news Web sites, a TOC block containing categorized news is usually similar in structure to that of block 1. Block 2 is also an informative TOC containing the abstracts of news and anchors linking to the news articles. Such useful evidences are more prominent in pages of the *systematic* Web sites in which ISs are usually generated automatically and dynamically by an iterative program from predefined templates.

In this paper, we extract and use knowledge from the tagging tree structure of a Web page and apply the Information Theory to mine the IS. Considering the structure information and context in these nodes together, we are able to understand and extract the meaning of information contained in Web pages more clearly and precisely. Specifically, we propose in the paper an IS mining

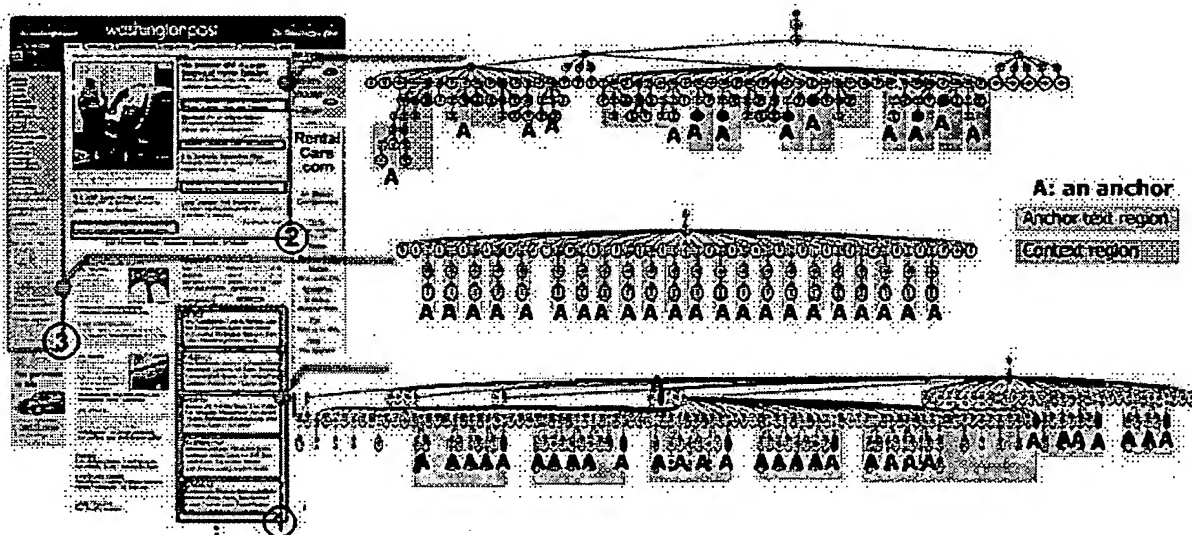


Fig. 2. A sample news page from WashingtonPost and the tree structures of informative blocks.

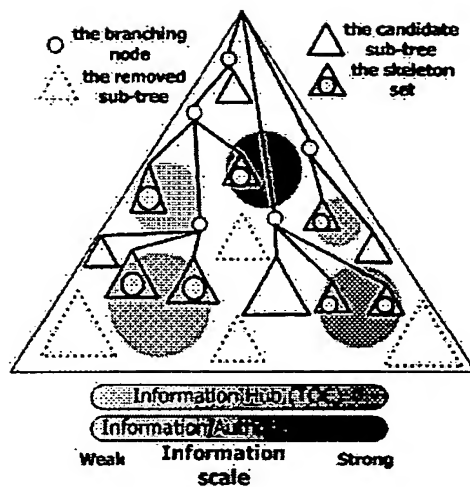


Fig. 3. The tree splitting and selection in WISDOM.

system called *WISDOM*, standing for *Web Intrapage Informative Structure Mining based on the Document Object Model* to automatically extract and recognize the *IS* of each page in a Web site.

The main mining flow in *WISDOM* first uses the Information Theory to evaluate the information amount contained in each DOM tree node and then constructs the *IS* by applying the specific searching, filtering, and merging methods. The searching step finds *IB* candidates and its principle is based on the observation that the root node of an *IB* uniformly spreads its information around its children nodes in most cases. In view of this, *WISDOM* first splits the original DOM tree into several small and nonoverlapped subtrees as shown in Fig. 3 and selects some of them as the candidate subtrees, in accordance with the assigned threshold of the structure information. The threshold is applied for the judgment on the uniformity of information distribution. Some uninformative subtrees are removed in the step. Our system then applies a top-down *IB* searching algorithm to find the top-*k* most informative blocks and the corresponding filtering criteria to select a set of candidate *IS*s called the skeleton set. The skeleton set can be considered the core subtrees of the *IS* shown in the color shaded regions in Fig. 3. The *IS* is built by expanding the skeleton set using the proposed merging methods. The merging method works in a bottom-up manner to link the qualified sibling nodes in the skeleton set and other informative nodes.

The remainder of this paper is organized as follows: In Section 2, we describe related work. *WISDOM* is described in Section 3. In Section 4, we evaluate the performance of *WISDOM* by testing it on several real news Web sites, university, and commercial Web sites. The Section 5 gives our conclusion.

## 2 RELATED WORK

Many works have been proposed that aim to extract the information of a page. Works on *wrappers* [9], [19], [22] provide learning mechanisms to mine the extraction rules of documents. The WebKB project in [6] automatically creates a computer understandable knowledge base from the textual content and hyperlink connectivity of Web pages. The work

describes an inductive logic programming algorithm for learning wrappers and develops a trainable information extraction system. Works in [1], [15], [20] provide auxiliary systems to aid in the information extraction from semistructured documents. The clipping method proposed in [14] is based on a supervised learning to provide a practical tool to cut the news articles. However, they need either a premarked training set or a considerable amount of human involvement to perform information extraction. When we consider the whole World Wide Web as our problem domain, building a useful training set to represent the diversity of Web content and structure is very hard.

In a systematic Web site, *IB*s are usually generated by a loop program; the entities in blocks are therefore similar to one another in view of their tag patterns and information they carry. In Fig. 2, it can be seen that the tag patterns of micro blocks (the shaded regions) in *IB*s 1 and 2 look very similar to one another. Therefore, frequent substructure mining is a candidate solution for automatic extraction of *IB*s. The topic of mining frequent substructure on the DOM trees of semistructure pages has recently been studied in [2], [10], [24] where the frequent subtree was extracted by respective pattern mining and noise node concealment methods, such as the wildcard mechanism in [10] and node-skip and edge-skip pruning in [2]. Works also use the tree pattern mining to extract metadata information in Web pages [13], [28]. However, semantic information in mined blocks with the same tree structure may be different from one to another. We need other information measurement methods to filter out redundant information blocks from those blocks with a similar tree structure. Moreover, some *IB*s like article blocks are laid out with the unique structures and are indeed difficult to extract by the frequent structure mining.

Some techniques proposed in [12], [29] use the semantics and relationships of tags to extract the record boundaries of Web pages. Several heuristic rules of tag characteristics, such as the highest count-tags (HT), identifiable "separator" tag (IT) and repeating tag pattern (RP), are proposed in [12] and are applied to extract record boundaries on several ".com" Web sites. Research in [29] also categorized tags into several groups according to their tagging functionalities and discovered the major schemas between them to translate HTML documents to XML documents in a semantic view.

Research in [7] extends the definition of a hub by dividing a hub page into several fine-grained hub blocks with different hub values. This is accomplished by calculating and integrating the hub values of each anchor node in the DOM tree of a page. Entropy analysis proposed in [21] discriminates the informative authorities of pages by dividing a page into several authority blocks with different authority values weighted by the information of each block.

There are also works on mining informative structure [16], [21], which are different from our work in that they mainly deal with mining blocks delimited by <TABLE> tags. In contrast, we mine fine-grained blocks using the DOM tree. It is worth mentioning that in the problem of mining the fine-grained *IB*s in a page, a straightforward approach would be to divide the page into several unit blocks that contain only one tag and then to merge neighboring blocks that contain information together. This naive method, however, does not work well for real-world Web pages in our opinion because: 1) When an *IB* is divided into small blocks with the one-tag granularity,



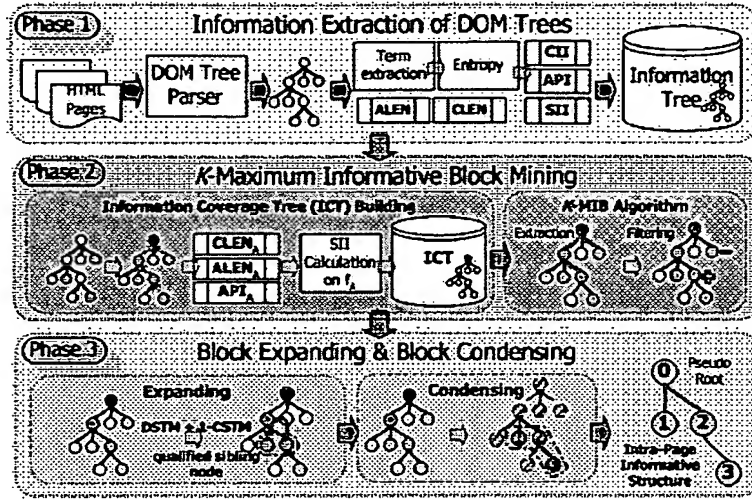


Fig. 4. WISDOM system flow.

the information contained is also divided into many small pieces which are difficult to discriminate from noises and redundant information and 2) there is no obvious method to merge such small blocks to form meaningful and integrated IBs. We therefore propose a top-down mining instead of bottom-up algorithm to extract fine-grained IBs.

### 3 WISDOM: A DOM-BASED MINING SYSTEM

WISDOM automatically extracts and recognizes ISs of each page in a Web site according to the knowledge in the tree structures of pages. As shown in Fig. 4, WISDOM consists of three phases: 1) information extraction from DOM trees, 2)  $k$ -maximum informative block mining, and 3) block expansion and condensation. In the first phase, we extract useful features from the information of the original DOM tree. These features can be classified into two types of information: node information and structure information. In the second phase, we aggregate the node information to build the Information Coverage Tree (ICT). According to the ICT, we devise a greedy algorithm, i.e.,  $k$ -maximum informative block mining algorithm ( $k$ -MIB), to extract subtrees that contain richer information. The extracted subtrees are either better information hubs or better information authorities, depending on the criteria employed in the greedy algorithm. They form the skeleton set of the IS of a page. We then expand the skeleton set by assembling neighboring subtrees that contain similar features corresponding to the original skeleton subtrees. After condensing the expanded set by removing dummy nodes, the assembled forest (or tree), in essence the IS of a page, is constructed.

#### 3.1 Phase 1: Information Extraction from DOM Trees

In the beginning, we crawl pages of a Web site in a specific crawling depth. When a page is crawled, we first extract the tree structure of a page based on DOM. Note that some HTML pages are not well-conformed, e.g., missing the ending  $\langle /a \rangle$  tag for the  $\langle a \rangle$  tag. We use HTMLTidy<sup>2</sup> to fix syntax mistakes in source documents. In tree  $T$ , each node

represents a tag in the page and contains the tag name information, attributes in the tag statement, and its *innerText*, i.e., the context delimited by the tag. From the definition of DOM, the context of the *innerText* of node  $N$  includes all contexts of nodes in the subtree rooted by node  $N$ . We use  $T(N)$  to denote the subtree rooted by node  $N$ . The *innerText* of the root node in each page is the context of a page when all tags are removed. The text of a Web page can be classified into two types: 1) anchor texts and 2) contexts which are texts delimited by all other tags except  $\langle A \rangle$  tags. We use  $ALEN$  to represent the length of the anchor text of a node and  $CLEN$  to represent the length of the contexts. A list of symbols used in this paper is given in Table 1.

We then parse the *innerText* of the root node to extract meaningful terms. A term corresponds to a meaningful keyword or phrase. Applying stemming algorithms and removing stop words based on a stop-list, English keywords (terms) can be extracted in a systematic manner [26]. Extracting terms in oriental languages is more difficult because of the lack of separators in these languages. In our system, we use an algorithm to extract keywords from Chinese sentences based on a Chinese term base. This base was generated by our search engine<sup>3</sup> by collecting hot queries and excluding stop words. After extracting terms in all crawled pages, we calculate the entropy value of each term according to its term frequency. From Shannon's information entropy [25], the entropy of term  $term_i$  can be formulated as:

$$EN(term_i) = - \sum_{j=1}^n w_{ij} \log_n w_{ij}, \text{ where } w_{ij} > 0$$

and  $n = |D|$ ,  $D$  is the set of pages,

where  $w_{ij}$  is the value of normalized term frequency in the page set. In the experiments on real Web sites containing a huge amount of pages, it is not practical to recalculate entropy values directly when a new page is crawled. In WISDOM, we use an incremental entropy calculation

2. HTMLTidy is an HTML fixing tool developed by Dave Raggett from the W3C team, <http://www.w3.org/People/Raggett/tidy/>.

3. The searching service is a project sponsored by Yam, a commercial search engine in Taiwan (<http://www.yam.com/>).



TABLE 1  
The List of Symbols Used

Abbr.	Description.	Abbr.	Description
ALEN	length of anchor text	ALEN <sub>A</sub>	aggregated ALEN
CLEN	length of contexts	CLEN <sub>A</sub>	aggregated CLEN
API	anchor precision index	API <sub>A</sub>	aggregated API
F	the set of tuple values, (ALEN, CLEN, API)	F <sub>A</sub>	the set of aggregated tuple values, (ALEN <sub>A</sub> , CLEN <sub>A</sub> , API <sub>A</sub> )
T	a DOM tree	ICT	tree T with the aggregated set F <sub>A</sub>
N	a node in the tree	T(N)	sub-tree rooted by N
innerText	contexts contained in T(N)	TLEN <sub>A</sub>	InnerText length
CII	content information index	SII	structure information index
ST	SII threshold	TC	type constraint
DSTM	direct sibling tree merging	k-CSTM	the k-th collateral sibling tree merging

process in the real Web site analysis. In the incremental calculation process, the new entropy value  $E_{k+1}(f_j)$  is calculated only by the previous entropy value  $E_k(f_j)$ , total term frequency  $TF_{j,k}$ , and the new term frequency  $tf_{(k+1)j}$  of the new included page for term  $f_j$ . The incremental calculation can be described as

$$E_{k+1}(f_j) = \begin{cases} \frac{E_k(f_j)}{\log_k(k+1)}, & \text{when } tf_{(k+1)j} = 0 \\ \Phi(E_k(f_j), TF_{j,k}, tf_{(k+1)j}), & \text{otherwise.} \end{cases}$$

The proof of the correctness of the process is given in Appendix A (which can be found on the Computer Society Digital Library at <http://computer.org/tkde/archives.htm>).

We define the weight of a term  $T_j$  as  $W(T_j) = 1 - EN(T_j)$  to represent the importance of the term. The reason behind applying entropy calculation is that terms distributed in more pages in a Web site usually carry less information to users. In contrast, those appearing in fewer pages carry more information of interest. The weight of a term is similar to its inverse document frequency, IDF [3], which is defined as  $\log_n \frac{n}{df_j}$ , where  $df_j$  is the document frequency of  $T_j$ . IDF is usually applied to represent the discriminability of a term in a set of documents. According to the definition, we can conclude following relationships between  $W(T_j)$  and  $IDF_j$ : 1) If  $T_j$  is uniformly distributed among some pages,  $W(T_j) = IDF_j$ . 2) If  $T_j$  is not uniformly distributed among the same pages in Item 1, then  $W(T_j) > IDF_j$  and the more skewed the distribution of  $T_j$  is, the larger  $W(T_j)$  is. The two relationships are proven in [17] and we include the detail of proofs in Appendix A which can be found on the Computer Society Digital Library at <http://computer.org/tkde/archives.htm>. Benefiting from these two relationships, the weight of a term attained from the entropy value is more representative for the importance of a term than from IDF. We use the example illustrated in Fig. 5 to explain these relationships. In this figure, Term<sub>A</sub> is uniformly distributed among Page 1 to Page 3 and Term<sub>B</sub> has the same term frequency and the document count with Term<sub>A</sub>, but most Term<sub>B</sub>s are located at Page 3. These two relationships are conformed by the following calculations:

$$\begin{cases} W(Term_A) = 1 - EN(Term_A) = 1 - 3 \cdot \frac{1}{3} \log_3 \frac{1}{3} = 0.207519 \\ = IDF(Term_A) = \log_3 \frac{4}{3} = 0.207519 \quad (1) \\ W(Term_B) = 1 - EN(Term_B) = 1 - 2 \cdot \frac{1}{6} \log_3 \frac{1}{6} - \frac{1}{6} \log_3 \frac{1}{6} \\ = 0.374185 > IDF(Term_B) = \log_3 \frac{4}{3} = 0.207519 \quad (2). \end{cases}$$

According to the extracted information, we calculate three extended features to gain more implicit information from the tree, namely, 1) the content information index (CII) which indicates the amount of information contained in the block, 2) the anchor precision index (API) which represents the similarity between the anchor-text and the linked document, and 3) the structure information index (SII) which indicates the distribution of children's feature values of one node in the DOM tree. Each node in DOM tree  $T$  contains the tuple values of the feature set  $F = \{ALEN, CLEN, API\}$ . In the following sections, we will describe their respective calculations.

### 3.1.1 Content Information Indices (CII)

When entropy values of terms are calculated, we average the weight values of terms in an *innerText* of node  $N$  to get the content information index of  $N$ , i.e.,

$$CII(N) = \frac{\sum_{j=1}^k W(term_j)}{k},$$

where  $\forall_{j=1 \sim k} term_j$  in *innerText* of  $N$ .

The CII value of node  $N$  represents the amount of information carried in a subtree rooted by  $N$ . The works in [16], [21] have shown that the entropy value corresponds to the recognition of the context parts of article pages. Consider the CII distribution of an article page in Fig. 6a. Note that the DOM tree is built by a depth-first traversal. The node ID of each node in the tree is generated according to the traversal order. Nodes with close node IDs are

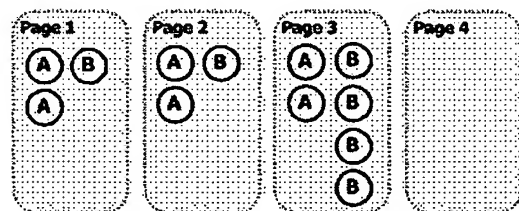


Fig. 5. An example of different term distributions.

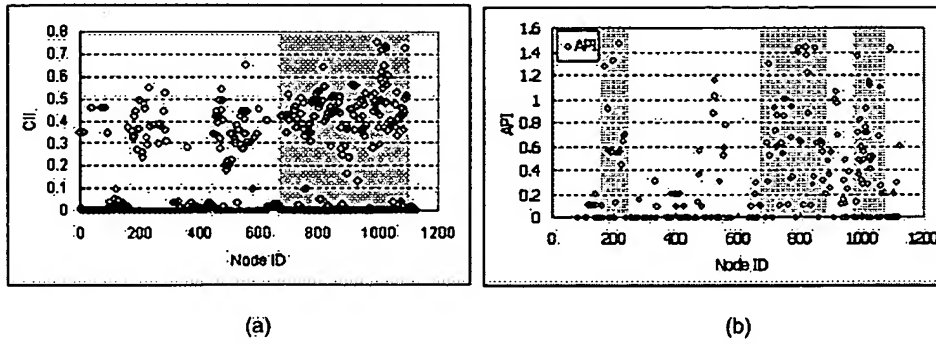


Fig. 6. The distributions of CII and API for two sample Web pages (shaded regions are IBs identified from the answer sets).

adjacent to each other in the physical layout of the page. The shaded region is an *IB* of the page, which is identified manually. Observe that 1) nodes with higher *CII* values in an *IB* are more than others and 2) nodes in an *IB* form a clear cluster in the *CII* distribution graph.

### 3.1.2 Anchor Precision Indices (API)

When browsing the Web, people use anchors to get information they want according to the semantics of anchors. The semantics of an anchor can be represented by the anchor text, text surrounding the anchor, the image, or other dynamic representations generated by scripts. The semantics of an anchor is expected to be relative to the page it links. Such relevance is, however, weak in some cases. We therefore define the value of the anchor precision index to indicate the correlation of the anchor and its linking page. We use the anchor text and the bounded text surrounding the anchor to evaluate the value of *API*. The correlation index *API* is defined as:

$$API(N) = \sum_{j=1}^m \frac{1}{EN(term_j)},$$

where  $term_j$  is the term concurrently appearing in both the anchor text of  $N$  and the linked page and  $m$  is the number of matched terms.

The calculation of *API* stems from the similarity analysis between documents using the vector space model. We extend the model by using the inverse values of entropy to set the weights of terms. If the information amount in those matched terms is larger, we get a larger *API* value that indicates that the anchor carries more precise information. The usage of the inverse of entropy values in the *API* formulation is to emphasize and amplify the effect of matched terms. Moreover, the value of *API* is not normalized by the matched count because we want to show that the longer informative anchor text leads to more information. Note that  $EN(term_i)$  is always larger than 0 because  $term_i$  appears in at least two documents.

Consider the *API* distribution of a TOC page in Fig. 6b. It shows that the number of nodes with larger *API* values in the shaded region, i.e., regions of marked TOC blocks, are more than others on average. Anchors in the menu block have small *API* values because the anchor texts of these anchors are short and the entropies of terms they contain are almost one.

### 3.1.3 Structure Information Indices (SII)

The index *SII* of a node is calculated according to the distribution of the feature values of the node's children. However, some HTML tags either correspond to information that is not extractable or provide no useful information. Such tags, such as the comment tag `<!-->`, the new line tag `<br>`, and the script program tag `<script>`, are called dummy tags and are removed from the following calculation of *SII*. We define the notion  $f_i(N)$  as the value of feature  $f_i$  of node  $N$ , and  $children(N)$  as the set of all nondummy children of the node  $N$ . For a simple tree structure of node  $N$  with children  $n_0, n_1, \dots, n_{m-1}$ , we define the *SII* value of node  $N$  for feature  $f_i$  as:

$$SII(N, f_i) = - \sum_{j=0}^{m-1} w_{ij} \log_m w_{ij},$$

$$\text{where } w_{ij} = \frac{f_i(n_j)}{\sum_{k=0}^{m-1} f_i(n_k)}, \forall n_k \in children(N).$$

Note that  $f_i(N)$  is larger or equal to the sum of  $f_i(n_0), f_i(n_1), \dots, f_i(n_{m-1})$ . We apply entropy calculation here to represent the distribution of children's feature values of any node with more than one child. The value of *SII* indicates the degree that the feature values of the node are dispersed among its children. When the value of  $SII(N, f_i)$  is higher, the values of all children's  $f_i$  tend to be equal.

In a systematic Web site, most context and anchors of TOC blocks are generated automatically. The styles, appearances, and information carried of entities in such a block are always similar from one to another. This phenomenon makes the *SII* values of these features become larger ones for the root nodes of such blocks.

## 3.2 Phase 2: The *k*-Maximum Informative Block Mining

In this phase, we first build the information coverage tree for features extracted during the phase one to obtain corresponding aggregated feature values. The proposed *k*-MIB algorithm is then applied to extract and filter out the candidate *IB*s. In Section 3.2.1, we describe the construction of *ICT* and the aggregated features. Extracting and filtering processes of the proposed algorithm are described in Section 3.2.2.

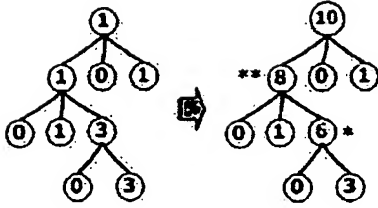


Fig. 7. An example of feature aggregation.

### 3.2.1 Information Coverage Tree Building

We define a tree with bottom-up aggregated features as an information coverage tree (abbreviated as ICT). In an ICT, any feature in the aggregated feature set  $F_A$  is obtained from the corresponding features in set  $F$ . Each node in an ICT contains all feature information of nodes in the subtree rooted by this node. The feature aggregation is a bottom-up process from the leaf nodes to the root node. The process of level  $k$  of the tree is shown below:

$$f_{Ai}(N) = f_i(N) + \sum f_{Ai}(n_j),$$

$$\forall n_j \in \text{children}(N) \text{ and } \text{level}(N) = k.$$

We aggregate features from the lowest level of the tree to the level one. The complexity of the process is  $O(|N|)$ . Fig. 7 shows an example aggregation process where the node marked by \* is labeled with  $6 = 3 + (3 + 0)$  and the one marked by \*\* is labeled with  $8 = 1 + (1 + 6)$ .

The aggregated features in ICT for each node  $N$  are subject to the constraint where  $f_{Ai}(n_j)$  is the aggregated value of feature  $f_i$  of node  $n_j$ :

$$f_{Ai}(N) \geq \sum_{j=0}^{m-1} f_{Ai}(n_j), \forall n_j \in \text{children}(N).$$

The length of *innerText* of each node is a typical aggregated feature because the *innerText* of a parent node contains all the *innerText* of its child nodes. We use  $TLEN_A$  to represent the length of *innerText*. In WISDOM, we also aggregate node information  $ALEN$  and  $API$  to get the corresponding aggregated features, denoted by  $ALEN_A$  and  $API_A$ . Note that  $TLEN_A(N)$  is composed of the length of contexts in  $T(N)$ , i.e.,  $CLEN_A(N)$ , and the length of anchor texts in  $T(N)$ , i.e.,  $ALEN_A(N)$ . The value of  $TLEN_A$  is thus equal to  $CLEN_A + ALEN_A$ . We then apply the SII calculation on these three aggregated features to get corresponding structure information of aggregated features for each node.

### 3.2.2 Block Extracting and Block Filtering

The proposed maximum informative block mining algorithm  $MIB(k, f_A, ST)$  is a greedy and top-down tree traversal process. For input value  $k$ , the algorithm outputs at most  $k$  IBs, i.e., TOC blocks or article blocks. The aim of the algorithm is to find the top- $k$  nodes with maximal aggregated feature  $f_A$  values under the given SII constraint, i.e., *SII Threshold* ( $ST$ ). When the value of  $ST$  is larger, the structure constraint is tighter and the children of each extracted node in the resulting candidate set will have more similar values of aggregated features in accordance with the

```

Algorithm MIB ( $k, f_A, ST$ ) begin
/* Cheap is a sorted stack */
1: InfoBlock = 0
2: Push root node into Cheap( $f_A$ )
3: While (InfoBlock < k and Cheap is not empty) begin
4:   Pop Node N with max( $f_A$ ) from Cheap( $f_A$ )
5:   IF ( $SII(N, f_A) > ST$  or N is a leaf) then
6:     find = true
7:     If (N matches the type constrain) then
8:       insert N into CandidateSet
9:       InfoBlock = InfoBlock + 1
10:    end if
11:   else
12:     push children(N) into Cheap( $f_A$ )
13:   end if
14: end
End

```

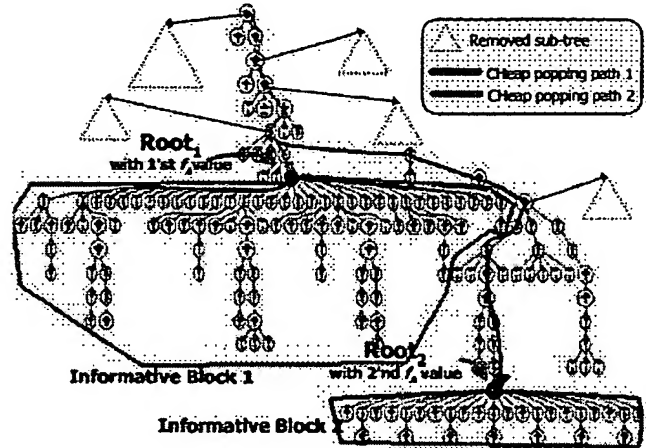


Fig. 8. An example of  $k$ -maximum informative block mining on the tree of a TOC page,  $k = 2$ .

definition of SII. The searching path of the algorithm is shown, for example, in Fig. 8. The original tree is extracted from a real TOC page by eliminating those subtrees removed by MIB.

When extracting the top- $k$  candidate nodes, we apply type constraints to eliminate pseudoinformative nodes. Type constraints (TC) are dependent on the type of blocks described as:

$$\begin{cases} \text{if type} = \text{"Article,"} & CII(N) \geq 1 - TC_{\text{article}} \\ \text{if type} = \text{"TOC,"} & \frac{API_A(N)}{\# \text{anchors in } T(N)} \geq TC_{\text{TOC}}, \end{cases}$$

$$\forall N \in \text{CandidateSet}.$$

Type constraints are motivated from heuristic observations that 1) article blocks contain informative context and, hence, their entropy values must be bounded and 2) TOC blocks contain highly semantic relevant anchors linking to information authorities and the average  $API$  value should be more than others in blocks with redundant and irrelevant anchors. These heuristic constraints are useful in removing pseudoinformative blocks.

Due to the tree traversal characteristic of the MIB algorithm, each node in the filtered candidate set is not an ancestor of any other nodes. The subtrees rooted by these nodes are therefore isolated and nonoverlapped. The set of these selected subtrees is called the *skeleton* of the IS of a page and the root nodes of these subtrees *skeleton nodes*.

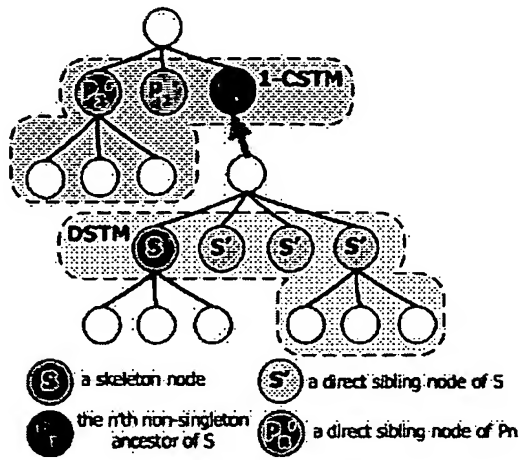


Fig. 9. Two sibling tree merging methods.

### 3.3 Phase 3: Block Expanding and Condensing

When investigating on the skeleton set, we find the selected skeleton nodes are often the subtrees of the IBs. The observation is more obvious when the structure threshold (ST) is larger and tighter. This is because the selected subtree is smaller when ST becomes larger in the  $k$ -MIB process. According to the skeleton structure, we therefore apply two sibling tree merging methods, i.e., direct sibling tree merging (DSTM) and collateral sibling tree merging (CSTM), to expand the skeleton set. The DSTM method merges the subtree rooted by qualified sibling nodes of each skeleton node  $S$  as shown in Fig. 9. Note that node  $S'$  may be one of skeleton nodes. Any qualified sibling node  $S'$  needs to match the type constraints and  $f_A(S')$  must also be smaller than  $f_A(S)$ . We do not need to merge sibling nodes with the larger  $f_A$  values because they are checked in the previous searching paths of the  $k$ -MIB algorithm and have been either selected into the skeleton set or removed from the IS. After DSTM, we then select the nonsingleton ancestors, i.e.,  $P_1, P_2, \dots, P_n$ , of  $S$  for the process of CSTM. The  $i$ th nonsingleton ancestor  $P_i$  is the  $i$ th ancestor of  $S$  which has more than one nondummy sibling node. The dummy node is defined as a node whose value of  $f_A$  is zero, e.g., the node with  $CLEN_A = 0$  for the article block and the node with  $ALEN_A = 0$  for the TOC block. The method of  $k$ -CSTM is equal to applying DSTM on  $P_k$ . In WISDOM, we apply DSTM and 1-CSTM to proceed the default block expanding. For example, in Fig. 9, we merge subtrees rooted by three  $S'$  into the skeleton set in the DSTM process. In the 1-CSTM process, we first traverse the tree from the node  $S$  up to the root node to find the first nonsingleton ancestor  $P_1$  and we then apply DSTM on  $P_1$  to merge its qualified sibling nodes, i.e., two  $P_1'$  nodes.

The intention of DSTM is to merge small IBs surrounding the skeleton blocks together. The nonuniform distribution between  $f_A$  values of the skeleton node and corresponding sibling nodes leads to the node separation in the  $k$ -MIB phase. In our experiments, DSTM can merge the metadata blocks, i.e., the article title, date, reporter, etc., into the main body of the article news. They are all IBs, but the distribution of their context length is skewed.

The block condensing process removes the subtrees rooted by nodes that cannot match the type constraints from

the expanding trees as these dummy subtrees are mainly tags for the page layout. This process is used to remove the uninformative subblocks from the merged trees obtained from the previous processes.

## 4 EXPERIMENTS AND RESULTS

In this section, we describe several experiments conducted on some real news Web sites in order to evaluate the performance of WISDOM. Data sets used and employed evaluation criteria are described in Section 4.1. We evaluate the performance of selection and filtering in the  $k$ -MIB algorithm in Section 4.2. The performance of block expanding and condensing is assessed in Section 4.3. Finally, Section 4.4 provides the overall performance evaluation of WISDOM.

### 4.1 Data Sets

We conduct our experiments on the data sets<sup>4</sup> used in [16]. In addition to these news Web sites, for evaluating WISDOM on other domains, a good example is the data set used in the WebKB project [8] and the data set used in the page segmentation research [9]. These data sets contain several university sites and commercial Web sites as described in Table 2. To assess WISDOM, we add two new answer sets, i.e., TOC blocks and article blocks. These blocks are extracted manually by news domain experts according to their experience in issuing real-world newspapers. We select most TOC pages and some candidate pages among all article pages with different tagging structures to mark. Unmarked TOC pages are pages which cannot be correctly parsed, or those containing many outside anchors linking to uncrawled pages, e.g., TOC pages in CNET and TTV. The latter case will cause the accuracy of API calculation to decrease suddenly and blur the evaluation results.

As shown in Table 2, the percentages of information coverage of the IS over the original page vary among data sets. Values are dependent on the styles and page layouts of news sites. The more redundant information added, the less information the IS carries.

To attain a quantitative evaluation, we employ two different evaluating methods to measure the values of precision and recall of article and TOC blocks. The TOC evaluation method is called significant node coverage (SNC). In SNC, we count the matched anchor nodes in subtrees rooted by nodes in the answer set and our output. For evaluating article blocks, we calculate the ratio of the matched context contained in each subtree by length to indicate the performance. The method is called information coverage (IC). The selection is made because only the context and anchors in the IS need to be indexed and extracted for crawling. In our experiments, we use the rates of precision (P) and recall (R) to indicate the similarity of these two sets. We also use F-measure [3] which combines recall and precision in a single efficiency measure. The value is the harmonic mean of precision and recall, and is formulated as  $\frac{2 \cdot (R \cdot P)}{R + P}$ . With the example in Fig. 10, we show the evaluation results of four methods in Table 3. The answer sets are two sets of root nodes, i.e., the TOC answer set  $A_T = \{a_{T1}, a_{T2}, \dots, a_{Tn}\}$  and the article answer set

4. Pages of Web sites in data sets were crawled on 2001/12/27, 2002/4/11, and 2004/3/29. The data sets can be retrieved at our research site: <http://kp06.iis.sinica.edu.tw/isd/index.html>.

TABLE 2  
Data Sets and Their Informative Structure Distributions

Site Abbr.	URL	Total pages	TOC pages	Marked TOC pages	Marked article pages	Marked Toc blocks	Marked article blocks	Answer Coverage	
								TOC (S-NC)	Article (IC)
CDN	www.cdn.com.tw	261	25	22 <sup>*</sup>	60 <sup>#</sup>	38	63	46.30%	98.40%
CTIMES	news.chinatimes.com	3747	79	69	66	313	68	32.10%	82.50%
CNA	www.cna.com.tw	1400	33	29	50	106	50	21.90%	80.10%
CNET	taiwan.cnet.com	4331	78	38	37	84	86	17.50%	63.60%
CTS	www.cts.com.tw	1316	31	19	53	21	80	54.80%	52.10%
TVBS	www.tvbs.com.tw	740	13	12	50	25	50	73.70%	56.90%
TTV	www.ttv.com.tw	861	22	18	42	20	75	20.10%	54.50%
UDN	udnnews.com	4676	252	243	52	674	106	28.00%	67.80%
CORN	www.cs.cornell.edu	1346	N/A <sup>S</sup>	14	14	24	18	45.42%	80.80%
UTEX	www.cs.utexas.edu	2935	N/A	11	10	11	10	45.02%	84.90%
WASH	www.cs.washington.edu	1526	N/A	16	10	23	10	79.04%	69.98%
WISC	www.cs.wisc.edu	2973	N/A	10	15	11	15	41.73%	77.08%
ABOUT	compnetworking.about.com	498	N/A	10	10	11	48	18.65%	43.54%
ECNET	reviews.cnet.com	500	N/A	10	10	10	10	38.89%	57.72%
ESPN	sports.espn.go.com	494	N/A	8	10	16	12	28.32%	58.39%
MONET	www.mo.net	261	N/A	9	10	9	13	35.78%	54.91%
XML	www.xml.com	807	N/A	10	10	10	26	43.46%	80.89%

\*: Unmarked TOC pages are removed from the TOC answer set due to the error occurring when parsing their DOM trees.

#: Domain experts selected the article pages with different and distinctive tagging styles to be the article answer set.

S: We only select some TOC and Article pages containing different structures for performance evaluation. We do not find all answers in English Web sites.

$A_A = \{a_{A1}, a_{A2}, \dots, a_{An}\}$ . The extracted results are the set of TOC blocks  $W_T = \{w_{T1}, w_{T2}, \dots, w_{Tn}\}$  and the set of article blocks  $W_A = \{w_{A1}, w_{A2}, \dots, w_{An}\}$ .

We show the result of the incremental entropy calculation in Fig. 11. In the figure, the value of the Y-axis means the ratio of the resulting entropy and the final entropy value calculated from the whole page set. We can find that the ratio difference is smaller than 0.1 when the corresponding document count is larger than 200. Therefore, in the practical usage, WISDOM can achieve a stable performance when the crawled page set is smaller than the whole page set of a Web site.

#### 4.2 Evaluation of $k$ -MIB

After the  $ICT$  of a page is built, we have to determine the searching ( $f_A$ ) and branching ( $ST$ ) criteria before applying  $k$ -MIB to the  $ICT$ . These selection criteria of  $k$ -MIB will affect the performance of the algorithm. In Fig. 12, we first conduct experiments to show the effects of different selection criteria for TOC blocks. We select  $ALEN_A$  and

$API_A$  for the searching criteria and corresponding  $SII$  values for the branching criteria. The result in Fig. 12 shows that using  $SII(API_A)$  for the branching criterion outperforms the one using  $SII(ALEN_A)$  when the selection criterion is to use a threshold of being equal to or smaller than 0.8. This is because  $API_A$  contains more information for discriminating the informative and redundant links than the length of anchor texts does.

We then apply the  $k$ -MIB algorithm to the  $ICT$  with the parameter pair  $(k, f_A, ST)$ . We use different  $ST$  values to control the number and granularity of the  $IBs$ . When the  $ST$  value is larger, more tighter and smaller blocks will be induced as shown in Fig. 13. Note that the average size of  $IBs$  in CTIMES is about three times as others and is out of the boundary of Fig. 13. This is due to the existence of big TOC blocks with entries of all categories of news in CTIMES. The sizes of these blocks are about 800 tags (nodes).

In the second phase of WISDOM, type constraint filtering plays an important role to remove the false-positive nodes. The selection of  $TC_{TOC}$  and  $TC_{Article}$  is made as follows: The distributions of the average  $API$  values, i.e., the criterion of the TOC type constraint, of top- $k$   $IBs$  in UDN are shown in

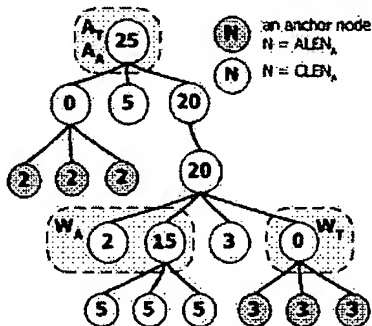


Fig. 10. A simple tree with an answer node and two results marked.

TABLE 3  
The Evaluating Calculation of the Example in Fig. 10

Method	$(A_T, W_T)$						$(A_A, W_A)$					
	AW	AO	WO	P	R	F	AW	AO	WO	P	R	F
SNC	3	3	0	0.50	0.50	0.67	5	5	0	1	0.50	0.67
IC	9	6	0	1	0.60	0.75	17	8	0	1	0.68	0.81

\*AW=the number of answer of the intersection of A and W  
 \*AO=the number of answer in A but not in W  
 \*WO=the number of answer in W but not in A  
 \*P = AW/(AW+WO), R=AW/(AW+AO)

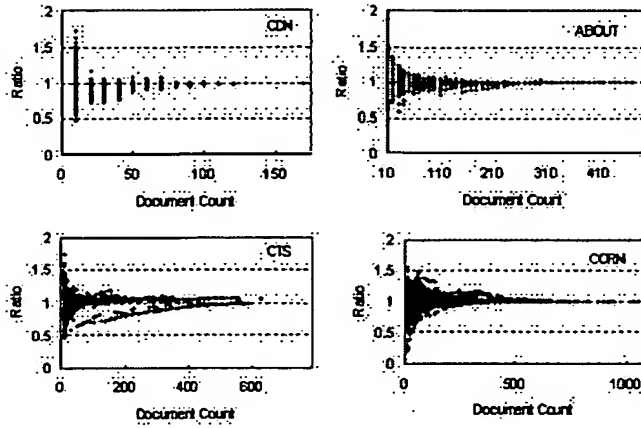


Fig. 11. Incremental entropy distribution for data sets CDN, CTS, ABOUT, and CORN.

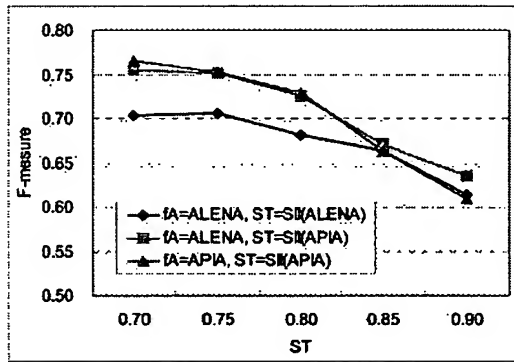
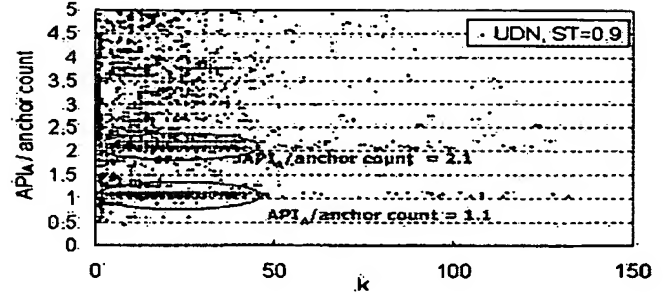

 Fig. 12. The effect of different criteria of  $k$ -mib for TOC blocks.

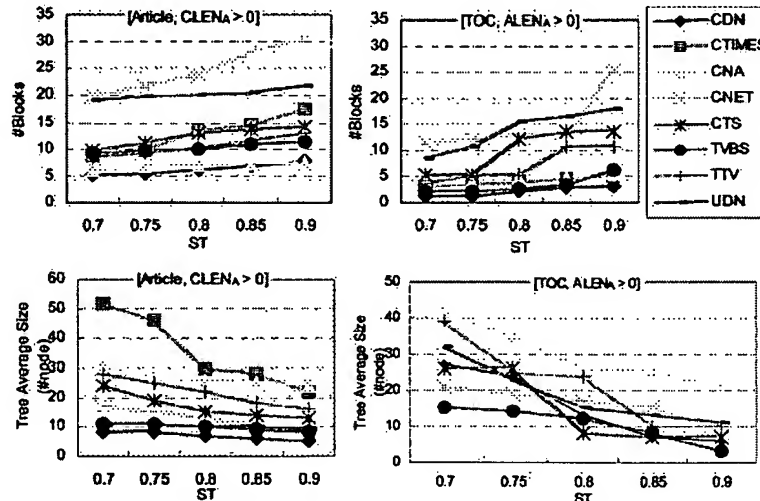
Fig. 14 where it can be seen that there are two obvious noise groups of values in this figure, i.e., 1.1 and 2.2, and they are reasonably chosen to be the  $TC_{TOC}$ . The selection of  $TC_{Article}$  is not so straightforward as the selection of  $TC_{TOC}$ . This is because when the size of  $IBs$  is divided into smaller ones, the number of extracted terms in each small block decreases, so does the accuracy of corresponding  $CII$ . Moreover, the index  $CII$  is not an aggregated value. We


 Fig. 14. The selection of the type constraint of TOC blocks [UDN,  $ST = 0.9$ ].

choose the uninformative link threshold described in [16], i.e., 0.8, to be  $TC_{Article}$ . Consequently, we use (1.25, 0.8) as default values for  $(TC_{TOC}, TC_{Article})$  for all data sets in WISDOM. The choice of  $TC_{TOC}$  value, 1.25, is simply motivated from the  $API$  formulation. We assume that each basic informative link contains one matching term with entropy 0.8, and its  $API$  value is 1.25 by the formulation. The value conforms to our observation on the real data shown in Fig. 14.

We show the average precision and recall values of  $k$ -MIB under the different selections of  $ST$  and  $k$  in Fig. 15. The results of TOC and article blocks both show the phenomena incurred by  $ST$ . When the value of  $ST$  increases, the sizes of split  $IBs$  decrease and the granularities of these blocks become finer. The selection of more fine-grained  $IBs$  increases the precision, but reduces the coverage of  $IBs$ , i.e., the recall. The same observation can be made in the same figure when the value of  $k$  becomes smaller.

From the results in Fig. 15, WISDOM is good at mining the informative article blocks rather than TOC blocks. First, there exists only one informative article block in most marked article pages. The information of an article page is thus more concentrated than information of a TOC page. This helps WISDOM discriminate informative article blocks easily. Second, noises affect and blur the  $API$  value. Using entropy to indicate the amount of information does not work well when few terms are extracted from the anchor


 Fig. 13. The average number and size of total  $IBs$  in a page selected by  $k$ -MIB without filtering.



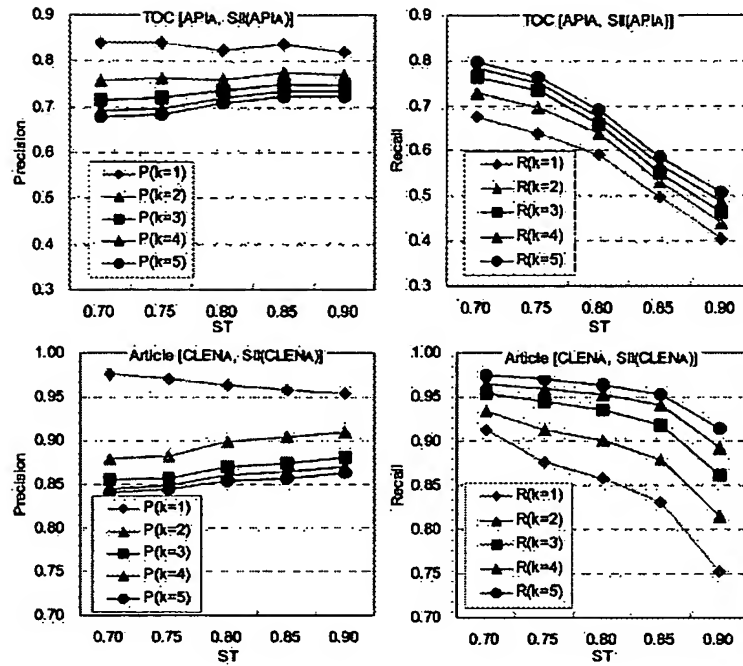


Fig. 15. The average values of precision and recall before phase 3 (caption: block type  $[f_A, SII]$ ).

text. The local menu effect mentioned in [16] can also decrease the entropy values of anchors in the menu blocks. The noise effects are prominent in CNET and UDN and the discriminability of the API value in these Web sites decreases suddenly. In CNET, more than 68 percent of all IBs have average API values of less than 2 if  $k \leq 5$ . Third, some informative TOC blocks mined by WISDOM are not "news" TOC blocks. These blocks are not selected in the answer set. The effects of applying different TCs are shown in Fig. 16. Filtering constraints can be used to remove pseudoinformative blocks.

### 4.3 Evaluation of Block Expanding and Condensing

Fig. 17 shows that the average improvement of different merging methods. The performances of experiments with different STs become similar after block expansion. This is because most blocks extracted by a high ST value are real IBs, though the sizes of these blocks are smaller than blocks extracted by a low ST value. The sizes of these smaller blocks can be expanded to the sizes of larger blocks by merging sibling subtrees which are also real IBs. Merging methods do not work well if a skeleton set contains many

pseudoinformative blocks, such as TOC blocks in CNET. Expansion of the skeleton set will incur more false-positive results. This is also the reason that the results with  $k = 1$  are better than those with  $k = 3$ .

### 4.4 Overall Performance

In Fig. 18, we use the system default setting, i.e.,  $k = 1$ ,  $ST = 0.8$ ,  $TC = (0.8, 1.25)$ , and merging methods DSTM and 1-CSTM, to show the overall performance of WISDOM on each data set. This figure shows that WISDOM is very good at the article blocks mining of all data sets and exhibits excellent performance on TOC blocks mining of CDN, CTIMES, CNA, CTS, and TVBS. The low values of precision and recall on CNET, TTV, and UDN are caused by the low accuracy of API values. Another low precision value on UDN is affected by the merging method 1-CSTM. Many pseudoinformative blocks are merged in the 1-CSTM step, even though WISDOM has reached the high recall rate after DSTM merging. The high average values of precision and recall also represent the robustness of WISDOM. We also compare WISDOM with two straightforward extracting methods in Fig. 19 to show the improvement. The method

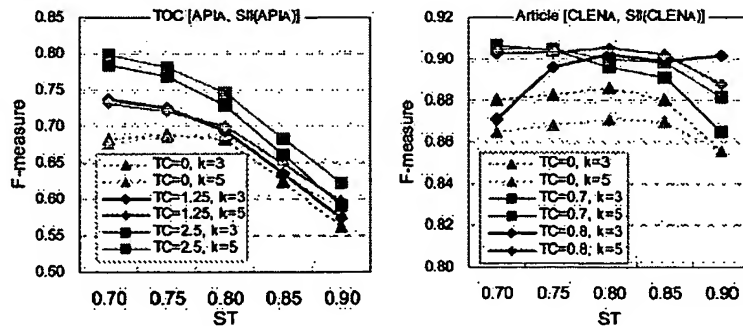


Fig. 16. The effects of different type constraints.

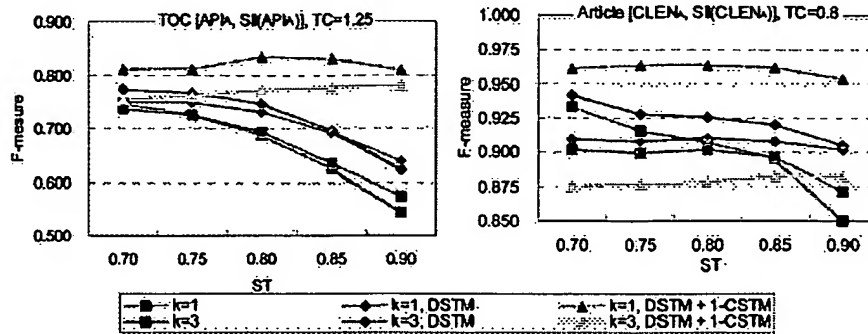


Fig. 17. The effects of DSTM and 1-CSTM.

M1 selects and merges leaf nodes of a DOM tree to a set of subtrees and is similar to the straightforward method described in Section 2. These merged leaf nodes must satisfy the same information constraint of WISDOM. The method M1 can be treated as simplified WISDOM that selects all leaf nodes into the skeleton set in the  $k$ -MIB phase. The method M2 works like M1 as well. The difference is that the method M2 uses the length constraint to filter the merged leaf nodes, i.e., the TLEN of a node must be larger than 5. Fig. 19 shows that WISDOM with the default setting leads these two methods and gives the prominent performance for the article pages.

The result in Fig. 20 shows the overall performance for English Web sites which consist of university and commercial domains by using the default setting same as in news Web sites. The performance for the article block extraction is also good as that in news Web sites. However, the performance of the TOC block extraction is worse than that in news Web sites. We found three reasons to cause the negative effect, which are 1) some informative anchors contains short anchor-text and common terms between anchors and linking pages are few. We cannot extract the anchor information in these cases and the feature APIs of these informative anchors therefore cannot be discriminated from redundant ones. 2) Some important terms are considered as stop-words and ignored, e.g., course-id in the university course pages. WISDOM ignores the numerical terms to reduce the noise effect caused by their high weights, which are obtained from the entropy calculation. However, course-id is an important clue for users to choice which course they feel interesting. This evidence shows that the stop-words selection must be dependent on the domain

characteristics, otherwise, some important words will be ignored. 3) Many commercial anchors are generated by the script language embedded in the pages. WISDOM cannot extract the anchor texts from these dynamic links.

To remedy these issues, we conduct an experiment to use different features instead of API to extract TOC blocks. The result in Fig. 21 shows the improvement when the feature ALEN and corresponding ST threshold are applied on some data sets in which API does not work well. This can be explained by that the TOC structure characteristics are retained when features ALEN and SII(ALEN) are applied and blurred when API is not correctly calculated or hard to be evaluated due to the lack of matched terms. API values are always zero even though corresponding ALEN values are larger than zero in these situations. We can evaluate the ratio of numbers of zero-API anchors over all meaningful anchors, i.e., their ALEN values are more than some threshold, in a page to be our criterion on the selection of appropriate features. The average ratios of data sets in Fig. 21 are obviously higher than others from our experimental observations.

## 5 CONCLUSION

We propose WISDOM to mine the ISs of a page. Given an entrance URL, WISDOM is able to crawl the site, parse pages into DOM trees, and analyze node and structure information in order to build information coverage trees. The system uses the Information Theory to split the DOM tree of a Web page into a set of IBs and uses the proposed searching ( $k$ -MIB), filtering, and merging (DSTM, 1-CSTM) methods to mine the IS of a page.

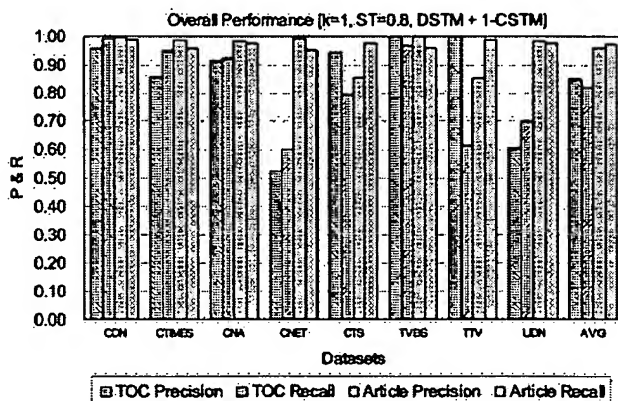


Fig. 18. Overall performance of WISDOM.

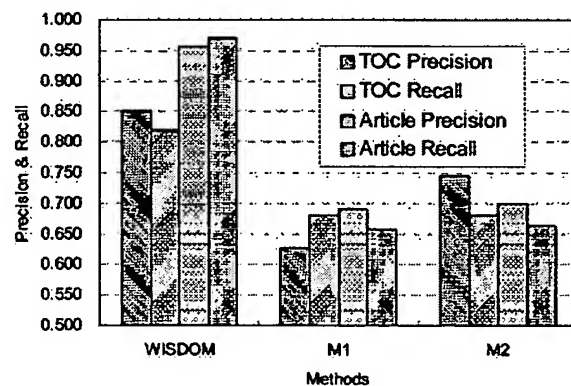


Fig. 19. Comparison of WISDOM to two straightforward methods.



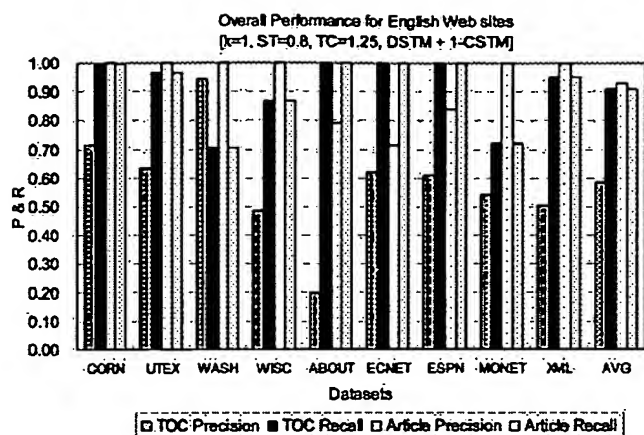


Fig. 20. Overall performance of WISDOM on other domain Web sites.

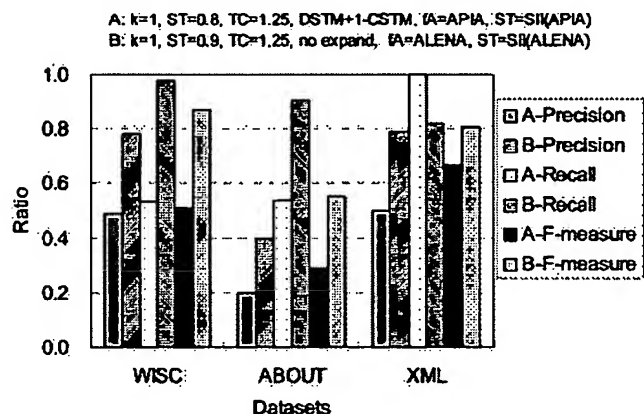


Fig. 21. Different feature selection for the TOC extraction.

For search engines, intermedia information agents, and crawlers, the IS mined by WISDOM is useful for indexing, extracting, and navigating significant information from a Web site. Experiments on several real news Web sites show high precision and recall rates attained by WISDOM which validates its practical applicability on news Web sites. We are integrating WISDOM into our news search engine (NSE) to help system managers speed up their work flow and reduce the labor of maintaining the site-dependent rule based extraction. For Web sites in other domains, even for nonsystematic Web sites, we are conducting some augmented feature to remedy the noise effects and improve the applicability.

## ACKNOWLEDGMENTS

The work was supported in part by the National Science Council of Taiwan, R.O.C., under Contracts NSC93-2752-E-002-006-PAE. H.-Y. Kao was with the Department of Electrical Engineering, National Taiwan University when this work was done.

## REFERENCES

[1] B. Adelberg, "NoDoSE—A Tool for Semi-Automatically Extracting Structured and Semistructured Data from Text Documents," *Proc. 1998 ACM SIGMOD Int'l Conf. Management of Data (SIGMOD)*, 1998.

[2] T. Asai, K. Abe, S. Kawasoe, H. Arimura, H. Sakamoto, and S. Arikawa, "Efficient Substructure Discovery from Large Semi-structured Data," *Proc. SIAM Int'l Conf. Data Mining (SDM)*, 2002.

[3] R. Baeza-Yates and B. Ribeiro-Neto, *Modern Information Retrieval*. Addison Wesley, 1999.

[4] Z. Bar-Yossef and S. Rajagopalan, "Template Detection via Data Mining and Its Applications," *Proc. 11th World Wide Web Conf. (WWW)*, 2002.

[5] A. Broder, S. Glassman, M. Manasse, and G. Zweig, "Syntactic Clustering of the Web," *Proc. Sixth World Wide Web Conf. (WWW)*, 1997.

[6] M. Craven, D. DiPasquo, D. Freitag, A. McCallum, T. Mitchell, K. Nigam, and S. Slattery, "Learning to Construct Knowledge Bases from the World Wide Web," *Artificial Intelligence*, vol. 118, nos. 1-2, pp. 69-113, 2000.

[7] S. Chakrabarti, "Integrating the Document Object Model with Hyperlinks for Enhanced Topic Distillation and Information Extraction," *Proc. 10th World Wide Web Conf. (WWW)*, 2001.

[8] Y. Chen, W.-Y. Ma, and H.-J. Zhang, "Detecting Web Page Structure for Adaptive Viewing on Small Form Factor Devices," *Proc. 12th World Wide Web Conf. (WWW)*, 2003.

[9] W. Cohen, "Recognizing Structure in Web Pages Using Similarity Queries," *Proc. Nat'l Conf. Artificial Intelligence (AAAI)*, 1999.

[10] G. Cong, L. Yi, B. Liu, and K. Wang, "Discovering Frequent Substructures from Hierarchical Semi-Structured Data," *Proc. SIAM Int'l Conf. Data Mining (SIAM SDM)*, 2002.

[11] R. Cooley and J. Srivastava, "Web Mining: Information and Pattern Discovery on the World Wide Web," *Proc. Ninth IEEE Int'l Conf. Tools with Artificial Intelligence (ICTAI)*, 1997.

[12] D.W. Embley, Y. Jiang, and Y.K. Ng, "Record-Boundary Discovery in Web Documents," *Proc. 1999 ACM SIGMOD Int'l Conf. Management of Data (SIGMOD)*, 1999.

[13] K. Furukawa, T. Uchida, K. Yamada, T. Miyahara, T. Shoudai, and Y. Nakamura, "Extracting Characteristic Structures among Words in Semistructured Documents," *Proc. Sixth Pacific-Asia Conf. Knowledge Discovery and Data Mining (PAKDD)*, 2002.

[14] H. Grundel, T. Naphtali, C. Wiech, J.-M. Gluba, M. Rohdenburg, and T. Scheffer, "Clipping and Analyzing News Using Machine Learning Techniques," *Proc. Int'l Conf. Discovery Science*, 2001.

[15] C.N. Hsu and M.T. Dung, "Generating Finite-State Transducers for Semi-Structured Data Extraction from the Web," *Information Systems*, vol. 23, no. 8, pp. 521-538, 1998.

[16] H.-Y. Kao, S.H. Lin, J.M. Ho, and M.-S. Chen, "Entropy-Based Link Analysis for Mining Web Informative Structures," *Proc. ACM 11th Int'l Conf. Information and Knowledge Management (CIKM)*, 2002.

[17] H.-Y. Kao, S.-H. Lin, J.-M. Ho, and M.-S. Chen, "Mining Web Information Structures and Contents Based on Entropy Analysis," *IEEE Trans. Knowledge and Data Eng.*, vol. 16, no. 1, Jan. 2004.

[18] J.M. Kleinberg, "Authoritative Sources in a Hyperlinked Environment," *Proc. ACM-SIAM Symp. Discrete Algorithms (SODA)*, 1998.

[19] N. Kushmerick, D. Weld, and R. Doorenbos, "Wrapper Induction for Information Extraction," *Proc. 15th Int'l Joint Conf. Artificial Intelligence (IJCAI)*, 1997.

[20] A. Laender, B. Ribeiro-Neto, A. Silva, and J. Teixeira, "A Brief Survey of Web Data Extraction Tools," *SIGMOD Record*, vol. 31, no. 2, June 2002.

[21] S.H. Lin and J.M. Ho, "Discovering Informative Content Blocks from Web Documents," *Proc. Eighth ACM Int'l Conf. Knowledge Discovery and Data Mining (SIGKDD)*, 2002.

[22] W.Y. Lin and W. Lam, "Learning to Extract Hierarchical Information from Semi-Structured Documents," *Proc. ACM Ninth Int'l Conf. Information and Knowledge Management (CIKM)*, 2000.

[23] X. Li, B. Liu, T.-H. Phang, and M. Hu, "Using Micro Information Units for Internet Search," *Proc. ACM 11th Int'l Conf. Information and Knowledge Management (CIKM)*, 2002.

[24] T. Miyahara, Y. Suzuki, T. Shoudai, T. Uchida, K. Takahashi, and H. Ueda, "Discovery of Frequent Tag Tree Patterns in Semistructured Web Documents," *Proc. Sixth Pacific-Asia Conf. Knowledge Discovery and Data Mining (PAKDD)*, 2002.

[25] C.E. Shannon, "A Mathematical Theory of Communication," *Bell System Technical J.*, vol. 27, pp. 398-403, 1948.

[26] G. Salton, *Automatic Text Processing: The Transformation, Analysis, and Retrieval of Information by Computer*. Addison Wesley, 1989.

[27] W3C DOM, Document Object Model (DOM), <http://www.w3.org/DOM/>, 2005.

- [28] K. Wang and H. Liu, "Discovering Structural Association of Semistructured Data," *IEEE Trans. Knowledge and Eng.*, vol. 12, no. 3, May/June 2000.
- [29] C. Yip, C. Gertz, and N. Sundaresan, "Reverse Engineering for Web Data: From Visual to Semantic Structures," *Proc. 19th IEEE Int'l Conf. Data Eng. (ICDE)*, 2002.



**Hung-Yu Kao** received the BS and MS degrees in the Department of Computer Science from National Tsing Hua University (NTHU), Hsinchu, Taiwan, in 1994 and 1996, respectively. In July 2003, he received the PhD degree from the Electrical Engineering Department, National Taiwan University, Taipei, Taiwan. From 1996 to 2003, he was a research assistant in the Institute of Information Science (IIS), Academia Sinica. He was a postdoctoral fellow of IIS from 2003 to 2004. Dr. Kao is currently an assistant professor of computer science and information engineering at National Cheng Kung University (NCKU). His research interests include Web information retrieval/extraction, knowledge management, data mining, bioinformatics, and content network.



**Jan-Ming Ho** received the BS degree in electrical engineering from National Cheng Kung University in 1978, the MS degree from the Institute of Electronics at National Chiao Tung University in 1980, and the PhD degree in electrical engineering and computer science from Northwestern University in 1989. He joined the Institute of Information Science, Academia Sinica, Taiwan, R.O.C., as an associate research fellow in 1989, and was promoted to research fellow in 1994. He visited the IBM T.J. Watson Research Center in the summer of 1987 and 1988, the Leonardo Fibonacci Institute for the Foundations of Computer Science, Italy, in the summer of 1992, and the Dagstuhl-Seminar on combinatorial methods for integrated circuit design, in October 1993. He is a member of the IEEE and the ACM. His research interests are targeted at the integration of theoretical and application-oriented research, including mobile computing, environment for management and presentation of digital archive, management, retrieval, and classification of Web documents, continuous video streaming and distribution, video conferencing, real-time operating systems with applications to continuous media systems, computational geometry, combinatorial optimization, VLSI design algorithms, and implementation and testing of VLSI algorithms on real designs. He is an associate editor of *IEEE Transaction on Multimedia*.



**Ming-Syan Chen** received the BS degree in electrical engineering from National Taiwan University, Taipei, Taiwan, and the MS and PhD degrees in computer, information, and control engineering from The University of Michigan, Ann Arbor, Michigan, in 1985 and 1988, respectively. Dr. Chen is currently a professor in the Electrical Engineering Department, a professor in the Computer Science and Information Engineering Department, National Taiwan University, Taipei, Taiwan, and the chairman of the Graduate Institute of Communication Engineering. He was a research staff member at IBM Thomas J. Watson Research Center, Yorktown Heights, New York, from 1988 to 1996. His research interests include database systems, data mining, mobile computing systems, and multimedia networking, and he has published more than 180 papers in his research areas. In addition to serving as program committee members in many conferences, Dr. Chen served as an associate editor of the *IEEE Transactions on Knowledge and Data Engineering* from 1997 to 2001, is currently on the editorial board of several journals, and is a distinguished visitor of IEEE Computer Society for Asia-Pacific from 1998 to 2000. He served as the program chair, program vice-chair, and general chair of several conferences. He was a keynote and tutorial speaker at several conferences, and also a guest coeditor for the *IEEE Transactions on Knowledge and Data Engineering* special issue on data mining in December 1996. He holds, or has applied for, 18 US patents and seven ROC patents in the areas of data mining, Web applications, interactive video playout, video server design, and concurrency and coherency control protocols. He is a recipient of the NSC (National Science Council) Distinguished Research Award and K.-T. Li Research Penetration Award for his research work, and also the Outstanding Innovation Award from IBM Corporate for his contribution to a major database product. He also received numerous awards for his research, teaching, inventions, and patent applications. Dr. Chen is a fellow of the IEEE and a member of the ACM.

► For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).


[CrossRef Search](#)
[Home](#) | [Login](#) | [Logout](#) | [Access Information](#) | [Alerts](#) |

Welcome United States Patent and Trademark Office

[BROWSE](#)
[SEARCH](#)
[IEEE XPLORE GUIDE](#)

You requested this document:

## » Key

**IEEE JNL** IEEE Journal or Magazine

**IEE JNL** IEE Journal or Magazine

**IEEE CNF** IEEE Conference Proceeding

**IEE CNF** IEE Conference Proceeding

**IEEE STD** IEEE Standard

## 1. Automatic web page classification in a dynamic and hierarchical way

Xiaogang Peng; Choi, B.;  
Data Mining, 2002. ICDM 2002. Proceedings. 2002 IEEE International Conference on  
9-12 Dec. 2002 Page(s):386 - 393

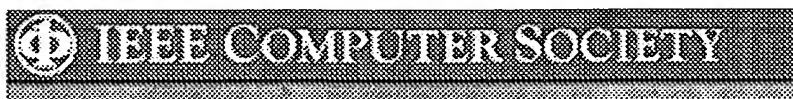
**Abstract:**

Automatic classification of web pages is an effective way to deal with the difficulty of re information from the Internet. Although there are many automatic classification algorithm that have been proposed, most of them ignore the conflict between the fixed number of the growing number of web pages going into the system. They also require searching t existing categories to make any classification. We propose a dynamic and hierarchical system that is capable of adding new categories as required, organizing the web pages structure, and classifying web pages by searching through only one path of the tree str results show that our proposed single-path search technique reduces the search comp increases the accuracy by 6% comparing to related algorithms. Our dynamic-category technique also achieves satisfying results on adding new categories into our system as

[Abstract](#) | [Full Text: PDF\(521 KB\)](#) **IEEE CNF**

 indexed by  
[Help](#) [Contact Us](#) [Privacy & ;](#)

© Copyright 2005 IEEE –

Search: 

Go

[Home](#)[Digital Library](#)[Site Map](#)[Store](#)[Help](#)[Contact Us](#)[Press Room](#)[Shopping Cart](#)[Login](#)

## digital library

### DIGITAL LIBRARY HOME

#### BROWSE BY TITLE

#### BROWSE BY SUBJECT

#### SEARCH

#### LIBRARY/INSTITUTION RESOURCES

#### RESOURCES

#### SUBSCRIPTION

#### ABOUT THE DIGITAL LIBRARY

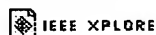
[Archive Page >>](#) [Table of Contents >>](#) [Abstract](#)

Second IEEE International Conference on Data Mining  
(ICDM'02) p. 386

### Automatic Web Page Classification in a Dynamic and Hierarchical Way

XIAOGANG PENG, Louisiana Tech University  
BEN CHOI, Louisiana Tech University

Full Article Text:



[BUY ARTICLE](#)

#### DOI Bookmark:

<http://doi.ieeeecomputersociety.org/10.1109/ICDM.2002.1183930>

#### Abstract

Automatic classification of web pages is an effective way to deal with the difficulty of retrieving information from the Internet. Although there are many automatic classification algorithms and systems that have been proposed, most of them ignore the conflict between the fixed number of categories and the growing number of web pages going into the system. They also require searching through all existing categories to make any classification. We propose a dynamic and hierarchical classification system that is capable of adding new categories as required, organizing the web pages into a tree structure, and classifying web pages by searching through only one path of the tree structure. Our test results show that our proposed single-path search technique reduces the search complexity and increases the accuracy by 6% comparing to related algorithms. Our dynamic-category expansion technique also achieves satisfying results on adding new categories into our system as required.

[Additional Information](#)

[Back to Top](#)

[Abstract](#)  
[Abstract](#)  
[Citation](#)

[Free ac](#)

☐ [Abstr](#)  
☐ [Selec](#)

[Electro](#)  
[in to](#)

☐ [Acce:](#)  
[text i](#)  
☐ [Dowr](#)  
[of PC](#)

[Subscr](#)

[Get a \](#)

**Citation:** XIAOGANG PENG, BEN CHOI. "Automatic Web Page Classification in a Dynamic and Hierarchical Way," *icdm*, p. 386, Second IEEE International Conference on Data Mining (ICDM'02), 2002.

---

Usage of this product signifies  
your acceptance of the Terms  
of Use.

This site and all contents  
(unless otherwise noted) are  
Copyright © 2002, IEEE, Inc.  
All rights reserved.

# Automatic Web Page Classification in a Dynamic and Hierarchical Way

XIAOGANG PENG & BEN CHOI  
Computer Science, College of Engineering and Science  
Louisiana Tech University, Ruston, LA 71272, USA  
[pro@BenChoi.org](mailto:pro@BenChoi.org)

## Abstract

*Automatic classification of web pages is an effective way to deal with the difficulty of retrieving information from the Internet. Although there are many automatic classification algorithms and systems that have been proposed, most of them ignore the conflict between the fixed number of categories and the growing number of web pages going into the system. They also require searching through all existing categories to make any classification. We propose a dynamic and hierarchical classification system that is capable of adding new categories as required, organizing the web pages into a tree structure, and classifying web pages by searching through only one path of the tree structure. Our test results show that our proposed single-path search technique reduces the search complexity and increases the accuracy by 6% comparing to related algorithms. Our dynamic-category expansion technique also achieves satisfying results on adding new categories into our system as required.*

## 1. Introduction

The World Wide Web is growing at a great speed but the documents in the Web do not form a logical organization and inevitably making the manipulation and retrieval difficult. The need for mechanisms to assist in locating relevant information becomes more and more urgent. One of the solutions to assist in retrieving documents on the Web is provided by classified directories [5]. However, current systems, such as Yahoo [30] still require human labor in doing the classification. Whether manual classification is able to keep up with the growth of the Web remains a question. First, manual classification is slow and costly since it relies on skilled

manpower. Second, the consistency of categorization is hard to maintain since different human experiences are involved. Finally, the task of defining the categories is difficult and subjective since new categories emerge continuously from many domains. Considering all these problems, the need of automatic classification becomes more and more important.

In this paper we present an automatic web page classification algorithm. The algorithm, unlike others, stresses the dynamic growing issue. It considers the hierarchical structure information for improving the classification accuracy. The core of the algorithm is a hierarchical classification technique that assigns a web page to a category. The number of web pages on the Web increases continuously in great speed and thus it is impossible for a fixed category set to provide accurate classification. To address this problem, we propose and implement a dynamic expanding technique.

### 1.1 Related Research

This paper relates to text learning and document classification. Text learning is a machine learning method on textual data that combines information retrieval techniques and is used as a tool to extract the content of textual data. A simple, yet limited, document representation (DR) is the "bag-of-words" text DR [12] [14]. Many experiments have been done to improve the performance of the DR. For example, Mladenic [18] extends the "bag-of-words" to the "bag-of-phrases" representation. Chan [4] also suggested that using phrases is a better choice than using single words. The goal of using phrases as features is to attempt to preserve the information left out by the "bag of words" methods. A document representation called "feature vector representation" uses a feature vector to capture the characteristics of the document by an "N-gram" feature

---

This research was supported in part by Center for Entrepreneurship and Information Technology (CEaIT), Louisiana Tech University, Grant iCSe 200123.

selection. An N-gram feature could be a word or a sequence of N words. A vector consists of a feature along with the occurrences of that feature within the document. Experiments show that N ranging from two to three is sufficient in most classification systems.

In information retrieval, TFIDF (Term Frequency-Inverse Document Frequency) classification algorithm is well studied [25]. Based on the document vector model, the distance between vectors is calculated by the cosine of the angle between them for the purpose of classification. Joachims [11] analyzed the TFIDF classifier in a probabilistic way based on the implicit assumption that the TFIDF classifier is as explicit as the Naïve Bayes classifier. By combining the probabilistic technique from statistic pattern recognition into the simple TFIDF classifier, Joachims proposed the PrTFIDF classifier with the formula:

$$P(d|c_j) = \sum_{w \in \{d, c_j\}} \frac{P(w|c_j) * P(c_j)}{\sum_{c \in C} P(w|c) * P(c)} * P(w|d) \quad (1.1).$$

Where  $c$  and  $c_j$  are categories taken from a category set  $C$ ,  $P(d|c_j)$  is the probability for a document  $d$  given a category  $c_j$ , and  $P(w|d)$  is the probability of a feature  $w$  given the document  $d$ .

The PrTFIDF classifier optimizes the parameter selection in TFIDF and reduces the error rate in five out of six reported experiments by 40%. Other more sophisticated classification algorithms and models were proposed including: multivariate regression models [8][26], nearest neighbor classifiers [31], Bayesian classifiers [15], decision tree [15], Support Vector Machines [7][10], voted classification [28]. Tree structures appear with all of these systems. Some proposed systems focus on classification algorithms to improve the accuracy of assigning testing documents to related catalogs [11], while others go even further by taking the classifier structure into account [12].

## 1.1 Organization

This paper is organized into the following sections: Section 2 describes our hierarchical classification module. Section 3 describes our dynamic expansion module. In Section 4 we use Yahoo structure as a test base and conduct several experiments to evaluate our system. Finally, in Section 5 we draw conclusions and provide future research.

## 2. Our Proposed Hierarchical Classification

In this section, we propose a new classification algorithm on a hierarchical structure. We first provide an overview of the algorithm then detail follows. The algorithm consists following stages: (1) generating category information tree, (2) hierarchical feature propagation, (3) feature selection on category information, and (4) single path traversal.

### 2.1 Generating Category Information Tree

For a web page classification system, the first step is to define the concept hierarchy using domain knowledge and to collect text data that corresponds to the concept hierarchy. The data is collected into an appropriate format for classification by a text-learning algorithm. The characteristics of each category are represented as a "bag of words" or a feature list based on the "well-grained text-learning method" [24]. Categories are arranged in a hierarchical tree structure. Each tree node represents one category. A child node of any given node represents a subcategory under the given node.

### 2.2 Hierarchical Feature Propagation

Many existing category structures are unbalanced hierarchical structures. In order to ensure the actual containment relation, the feature information of a category is propagated upward from leaf nodes to the root node of our classification tree. For example, Mladenic [22] studied Yahoo unbalanced structure and proposed an algorithm for featuring propagation. The algorithm takes care of the structure of the tree hierarchy. As proposed in [22], with a tree  $T$  rooted at node  $N$  having  $k$  sub-trees ( $SubTi$  ( $i=1...k$ )), we can calculate the probability for a feature  $w$  belonging to a category after propagation as follows:

$$P(w|T) = \sum_{i=1}^k P(w|SubTi) * P(SubTi|T) + P(w|N) * P(N|T) \quad (2.1)$$

Where  $P(w|T)$  is the propagated feature probability given tree  $T$  and  $P(w|N)$  is the probability of the feature  $w$  in the node  $N$  before the propagation, which can be calculated by dividing the particular term frequency by the total term frequency.  $P(SubTi|T)$  and  $P(N|T)$  are the weight factor of the sub-tree  $SubTi$  given tree  $T$  and the probability of current node given tree  $T$ .

In the generated classification tree (from stage 1), each node represents a category by a feature list. The propagated feature list is generated by adding the original feature list of the current node and all the propagated feature lists of the sub-categories and by assigning different weights. By propagating in this way, the feature list captures the characteristics of a sub-tree rooted in the current node rather than in an isolated category set.

### 2.3 Feature Selection on Category Information

It is clear that what really contributed in distinguish between categories are those unique features belonging to the categories. Because of the feature propagation, these unique features will be weighted and propagated upwards and become the features of the parent category. In this case, by tracing these unique features it is easy to locate the correct category. Due to the uniqueness of the features, there is only one path in the tree for reaching the features. This phenomenon provides a foundation for our single path classification algorithm.

Mladenic and Grobelnic [20] conducted similar experiments in feature selection on Yahoo category tree and proposed Odd Ratio measurement, but they used the complement of a node from the entire tree as negative examples to calculate the Odd Ratios.

The goal of our feature selection is to compare the features in a node to its sibling nodes and try to distinguish the unique features. In order to do this, we take advantage of the feature propagation and use the features of the parent node as negative examples to determine a ranking. After propagation, each propagated feature probability is actually the weighted sum of the same feature probability from the sub-trees and the current node itself. We proposed the following formula for determining the uniqueness ranking  $R_c(w)$  of a feature  $w$  of a category  $c$ .

$$R_c(w) = \frac{P(w|c) * P(SubT_c|T)}{P(w|parent)} \quad (2.3)$$

Where  $c$  is the current node and *parent* is the parent node of  $c$ .  $P(SubT_c|T)$  is the weight factor that assigned to node  $c$  when it is propagated to the parent. If a feature is unique in one node, it is the only source that can propagate to the parent feature list. In this case, we get the maximum value of the formula, which is 1. The unique features will be considered as key features that differentiate the node from its siblings and forms the basis for our single path traversal algorithm.

### 2.4 Single Path Traversal

The concept of text classification requires the use of a classifier to assign values to a document and to each catalog. Matching a document feature values to a category feature values, the catalog with a global maximum value is considered the correct place to hold the document. Most automatic classification researches concentrated on the global search algorithm. They treat all categories in a flat structure when trying to find the maximum. It follows that in order to find the category with the greatest value, it is necessary to compare all the categories.

In the tree structure that we have configured, we claim that traveling one path is sufficient to achieve this goal. If there are  $N$  categories in the tree, the complexity of searching all categories is  $\theta(N)$ , but by our single path algorithm it is merely  $\theta(\log(N))$ .

The idea of the single path traversal is to eliminate the impact of other branch in the tree. After feature propagation, the propagated feature list of the parent node is a scaled summation of the propagated feature lists of its children. Thus, by checking the parent node we can know the information of its descendents. In our tree structure, all the features are propagated upwards with the ranking function of formula 2.3, so each category has unique features of its own to differentiate from its siblings. These two steps make the single path traversal possible.

The first step of the single path traversal is to discriminate sibling nodes in each level and find a correct path for the incoming web page. In order to determine the discriminating probability for each node, we only consider the nodes at each level and apply the PrTFIDF formula 1.1 on the features with a ranking of 1, which indicate a unique feature. At each level, we chose the node with the maximum discriminating probability as the starting point for the next iteration. Recursively applying this rule creates a path from the root of the tree to one of the leaf nodes.

Then following this path we apply the PrTFIDF classifier again using all the features of the nodes belonging to the path, to get the actual probability for the page with categories within this path. By picking the node along the classification path with maximum actual probability value, we determine the candidate category for the page.

### 3. Dynamic Expansion and Updating

In this section, we describe our proposed new dynamic tree expansion algorithm. As more and more documents being put into a catalog set, the diversities of the



documents make the existing catalogs unable to guarantee classification accuracies. The problem of how to dynamically generate more sub-catalogs for the existing catalog set becomes evident. The highlight of our approach is as follow; details are provided in following subsections. Based on statistical results, we determine a set of criteria and check whether the criteria have been met for putting a page under the current category. If the criteria have not been met, we will create a new node for the page. As the results, the tree will be expanded by adding a new category. An updating algorithm is also introduced to incorporate the expansion information into the existing catalog set.

### 3.1 Dynamic Expansion

As the number of pages that need to be classified grows larger and because of their diversity, the original number of categories does not always fit the true content of the incoming pages. It is necessary for expanding the category tree to accommodate all the pages in order to yield accurate results. The criteria for creating new categories must be established. There are two types of expansion, deeper and wider, and two thresholds will be used to determine the criteria.

The two thresholds  $B$  and  $D$  are obtained in the following ways. We take a set of sample pages from a category and calculate the probability of these pages relating to its category. We denote the resulting values for each category as  $S_i$  for  $i$  from 1 to the number of all categories  $n$ . Then, we compute the normal distribution of all sample  $S_i$  for  $i = 1$  to  $n$ , and obtain the mean  $\mu$  and the standard deviation  $d$ . We set  $B = \mu - d$  and  $D = 2d$ .

We define the relation between a document and the category it should be placed as the "belongs to" relation. It is a basic assumption that the probability of a page given a category having the "belongs to" relation will have a maximum value. In expansion, we set a lower bound for this relation, threshold  $B$  will be used to ensure the maximum characteristic of the relation.

A deeper expansion happens when the maximum actual probability value (obtained from single-path search described in the last subsection) is smaller than threshold  $B$ . That is although the value is the maximum found, its corresponding category is not considered to be sufficiently suitable for the new page. A new sub-category under the category is then created to hold the new page.

A wider expansion occurs in the following situation: the probability of the page comparing to that the candidate category (the one with the maximum actual

probability value) is much bigger than the probability of the page and that the candidate's children categories. The difference of probability represents the distinction between the page and the candidate's children nodes. Ignoring this distinction will cause the inconsistency of the "belongs to" relation between parent and children. When updating the category information after a page is put into a category, those features, which contribute to the difference, will also be incorporated into the category. This makes the relation between the category and the existing children more and more distinct. By creating the new category, we put the distinction to the siblings other than making the relation between parent and children far apart. Thus, the newly added features will only heavily affect the new created category. When propagating the new features to the candidate category, the weight factor is used to reduce the effect of the new feature. This reduces the inconsistent effect.

Based on these two cases, the threshold  $B$  ("belongs to" threshold) and the threshold  $D$  (difference threshold) are used as criteria to determine when expansion is needed. For deeper expansion, the probability for the new page (document  $d$ ) given the candidate node is less than  $B$ , that is  $P(d|c) < B$ . In this case, we will create a new category. For wider expansion, the candidate node is not a leaf node of the tree. We need to check the probability of the page given each sub-categories of the candidate category  $P(d|Sub_i)$ . If the difference of  $P(d|c)$  and the maximum number of the  $P(d|Sub_i)$  is out of the range of threshold  $D$ , that is  $P(d|c) - \text{Max}(P(d|Sub_i)) > D$ , then the new page is considered to be substantially different from any of its sub-categories. In this case, we will create a new category.

### 3.2 Updating Category Information

When a page is assigned to a category, whether a page is just created or an existing one, the feature vector of the page will be contributed to the category information. It is necessary to update the category information feature list to maintain the concordance and the hierarchical structure and expansion of the vocabulary will inevitably change the characteristic of some of the categories.

After the page has been put into a category, the page is considered to be one part of the category, so it will contribute its own characteristics to the category. Both the page and category information are represented as feature vectors. Merging the page vectors into the category vectors is a solution for updating the category information. Since the page features are selected by a "well-grained text-learning method" [24] with those low frequency patterns removed, we assume that the extracted

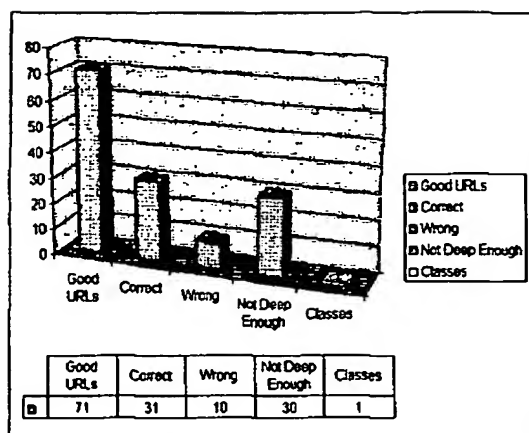


Figure 4.1 Results for accuracy of deeper expansion

features are considered to be relevant to represent the page. We use the following formula for updating:

$$P'(w|N) = \frac{|V| * P(w|N) + |V_{page}| * P(w|Page)}{|V| + |V_{page}|} \quad (3.1)$$

Where  $|V|$  represents the size of the category vector and the  $|V_{page}|$  is the page vector size.

After merging page vectors, the category information is changed, so the ranking and propagated feature probability should be recalculated. Since those changes happen just along the classification path, the updating takes place only along the single path.

#### 4. Experiments and Results

We design experiments to test two aspects of our system: accuracy in dynamic expansion and performances of the single path traversal. The root of our experiment is set to Yahoo /Science/Engineering/, one of the sub-categories of Yahoo's classification tree. We also chose the categories that are strictly following the levels of this root category; that is, we eliminated those categories that either go to another root category or do not follow the level structure. As we have noticed, many of the outgoing URLs in Yahoo are not accessible. We chose the Science/Engineering as the root because it is newly generated and covers 4068 outgoing URLs as advertised. Our expectation is that it will somewhat reduce the number of outdated outgoing URLs and provide enough testing examples for our system.

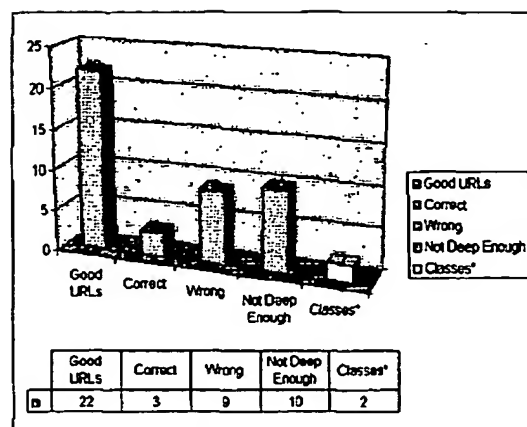


Figure 4.2 Experiment results of wider expansion

#### 4.1 Machine Learning Setting

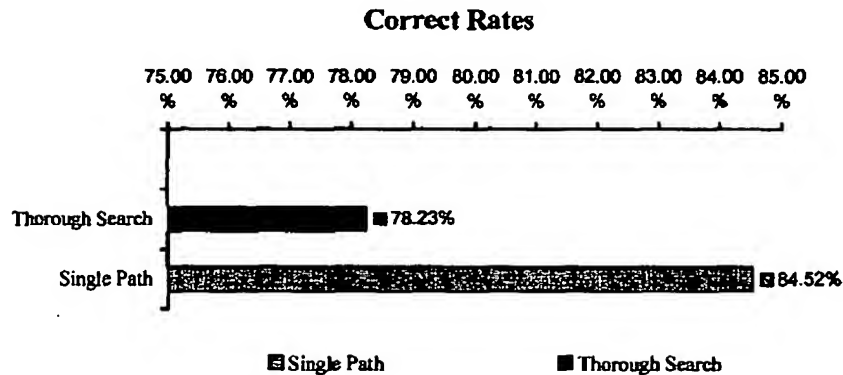
Considering processing time for testing purposes, we direct our program in getting the category information of three levels in the Yahoo "Science/Engineering" sub-tree. Labrou and Finin [13] compared several different ways in describing the category information and the web page information. By their experiments, they pointed out that, the best way of describing category was entry summaries and entry titles; correspondingly the best choice of web page description (called entry in their paper) was the entry summaries. Because of this, when generating the category information, we use summaries that are generated by man power and already there in Yahoo website, instead of using the actually website contents to generate the category information tree.

The testing examples of our system are actual website contents whose URLs are taken from three levels of the Yahoo sub-tree rooted in "Science/ Engineering." Some non-text format paged associated with the URLs cannot yet be classified, for example, jpeg, swf, and gif files. We only use the URLs whose pages having more than 70 features after our well grain 3-gram feature extraction, and these URLs are called "good" URLs.

The two global variables (thresholds) are generated at the machine-learning step using the statistic of the page-category probability. All the categories at the second and the third levels of "Science/ Engineering", except the "Science/ Engineering/organizations" category, are defined as existing categories. One third of the URLs that belong to those categories are taken to calculate the two thresholds based on the ways that we have described in section 3.1.

**Table 4.1 Results of comparing two algorithms**

Same classification results		
Correct results	Wrong results	
333	6	
Different classification results		
Correct results by Single path algorithm	Correct results by thorough search algorithm	Wrong results by both algorithm
58	30	37



**Figure 4.3 Correct rates of both algorithms**

## 4.2 Expansion Tests

Our expansion experiments are designed to test two kinds of accuracies: deeper test – testing the accuracy for deeper expansions, (see results in Figure 4.1), and wider test – testing the accuracy for wider expansions, (see results in Figure 4.2).

All of the URLs that have the “belong to” relation to the category “Science/Engineering/organization” are considered to be a “new” group and used to test the wider expansion. In this experiment, the number in “correct” column means how many pages are classified to an expanded category that is under the selected “Science/Engineering/organization”. If the page is classified to a category that contains the category where this page comes from, it is called “Not Deep Enough.” The column “Classes” keeps the number of the newly expanded category. Since all the URLs are taken from a same category, we are expecting the number to be 1. In deeper test, the testing URLs come from “Science/ Engineering/ Civil\_Engineering /Institutes/.” Similar for wider test, we

have same settings for the experiment results. Experiment results are provided on figure 4.1 and 4.2.

## 4.3 Testing our Single Path Algorithm

In order to test the accuracy of single path algorithm, we compare the actual classification results of the single path algorithm to the one that searches all categories. Using 33% of all the web pages in Yahoo “Science/Engineering” tree as testing web pages, we compare the two algorithms by the accuracy and the effectiveness, and the result is shown in Table 4.1 and Figure 4.3.

By the results, we can see that the two algorithms have the same result is more than 73% of all testing cases. Surprisingly, even when we have ignored most of the branches of the tree, the single path classification still has an accuracy of 84.26%, which even outperforms by more than 6% the accuracy that was achieved by the thorough search algorithm. From these results, we can see the advantage of our single-path search algorithms.

## 5. Conclusion

In this paper, we describe an approach that utilizes class hierarchies for improving text classification. Our single path classification algorithm, in the hierarchical classification module, reduces the computational expense compared to the thorough search algorithm that is used by most of the existing classification algorithms. By distinguishing the siblings, the algorithm recognizes a correct path containing the destination category. The algorithm is successful not only in saving computational resources but even improving the correct hits to a higher percentage. Our experiment also shows that because of the diversity of contents of the pages, the thorough search algorithm sometimes cannot tell the key information from the common information since it treats all the information equally. The single path algorithm avoids this problem by only considering unique features in discriminating siblings. Our experiments support that the single path algorithm is more competent both in time and in accuracy.

The expansion and updating modules emphasize the developing of a dynamic system. The expansion algorithm creates a new category when it detects that the page is not similar enough for the current node or might affect the containment relation between current node and its children. The updating algorithm keeps the tree database in a valid state. With expansion and updating, our database grows more diverse in order to proficiently categorize more incoming web pages.

Finally, since the Internet is growing in great speed but the arrangement of the web pages do not have a logical or semantic organization, to find a structured semantic Internet, we should find a standard organization for all Internet data. Since Internet resources can be considered as different kinds of information units, for grouping proposes, classification is perhaps the most appropriated way to organize them. Our system provides a starting point. Hierarchical structure and dynamic growing mechanism can be considered as bases of a structured Internet. We expect that in the near future, when all Internet resources can be classified, the whole Internet will be converted to a well-structured system based on certain classification standard. To achieve this goal much future research remains to be done.

## References

- [1] AltaVista, <http://altavista.digital.com>
- [2] A. Berker and V. Mittal, "OCELOT: a system for summarizing web pages", In Proceedings of SIGIR, 144-151, 2000
- [3] M. CatePazzani and D. Billsus, "Learning and Revising User Profiles: The Identification of Interesting Web Sites", *Machine Learning* 27, 313-331, 1997
- [4] Philip K. Chan, "A non-invasive learning approach to building web user profiles", KDD-99 Workshop on Web Usage Analysis and User Profiling, 1999
- [5] Chandra Chekuri, Michael H. Goldwasser, Prabhakar Raghavan and Eli Upfal, "WebSearch Using Automatic Classification", In Proceedings of the Sixth International World Wide Web Conference, 1997
- [6] P. Domingos and M. Pazzani, "On the optimality of the simple Bayesian classifier under zero-one loss", *Machine Learning* 29, 103-130, 1997
- [7] Susan Dumais, Hao Chen, "Hierarchical Classification of Web Content". Proceedings of SIGIR-00, 23rd ACM International Conference on Research and Development in Information Retrieval, 2000
- [8] N. Fuhr, S. Hartmann, G. Lustig, M. Schwantner, and K. Tzeras, "A rule-based multi-stage indexing system for large subject fields", Proceedings of RIAO'91, 06-623, 1991
- [9] HotBot, <http://www.hotbot.com>
- [10] Thorsten Joachims, "Text categorization with support vector machines: Learning with many relevant features", Proc. 10<sup>th</sup> European Conference on Machine Learning (ECML), Springer Verlag, 1998
- [11] Thorsten Joachims, "A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization", In International Conference on Machine Learning (ICML), 1997
- [12] D. Koller and M. Sahami, "Hierarchically classifying documents using very few words", Proceedings of the 14<sup>th</sup> international Conference on Machine Learning ECML98, 1998
- [13] Yannis Labrou and Tim Finin, "Yahoo! as an ontology – using Yahoo! Categories to Describe Document" In CIKM '99. Proceedings of the Eighth International Conference on Knowledge and Information Management, 180-187, ACM, 1999
- [14] K. Lang, "Newsweeder: Learning to filter news", In Proceedings of the 12th International Conference on Machine Learning, 331-339, 1995
- [15] D.D. Lewis, and M. Ringuette, "A comparison of two learning algorithms for text categorization", Third Annual Symposium on Document Analysis and Information Retrieval (SDAIR'94), 81-93, 1994
- [16] Lycos, <http://www.lycos.com>

- [17] Daniel Marcu, "From Discourse Structures to Text Summaries", Proceedings of the ACL/EACL-97 Workshop on Intelligent Scalable Text Summarization, 1997
- [18] Dunja Mladenic and Marko Grobelnik, "Word sequences as features in text-learning", In Proceedings of ERK-98, the Seventh Electro-technical and Computer Science Conference, 145-148, 1998
- [19] Dunja Mladenic, "Feature subset selection in text-learning", In Proceedings of the 10th European Conference on Machine Learning ECML98, 1998
- [20] Dunja Mladenic and Marko Grobelnik, "Feature selection for classification based on text hierarchy", Working Notes of Learning from Text and the Web, Conference on Automated Learning and Discovery, 1998
- [21] Dunja Mladenic, "Turning Yahoo! into an Automatic Web-Page Classifier", Proceedings of the 13th European Conference on Artificial Intelligence ECAI98, 473- 474, 1998
- [22] Dunja Mladenic, "Machine Learning on non-homogeneous, distributed text data", PhD thesis, University of Ljubljana, Slovenia, 1998
- [23] C. D. Paice, "Constructing Literature Abstracts by Computer: Techniques and Prospects", In Information Processing & Management, 26(1), 171-186, 1990.
- [24] Xiaogang Peng, "Automatic Web Page Classification in a Dynamic and Hierarchical Way", MS Thesis, Louisiana Tech University, 2002
- [25] G. Salton, and C. Buckley, "Term Weighting Approaches in Automatic Text Retrieval", Technical Report, COR-87-881, Department of Computer Science, Cornell University, November
- [26] H. Schutze, D. Hull, and O.J. Pedersen, "A comparison of classifiers and document representations for the routing problem", Proceedings of the 18<sup>th</sup> Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, 229-237, 1995
- [27] C.J. van Rijbergen, D.J. Harper, and M.F. Porter, "The selection of good search terms", Information Processing & Management, 17, 77-91, 1981
- [28] S.M. Weiss, C. Apte, F. Damerau, D.E. Johnson, F.J. Oles, T. Goets, and T. Hampp, "Maximizing text-mining performance", IEEE Intelligent Systems, 14(4), 63-69, 1999
- [29] Michael J. Witbrock and Vibhu O. Mittal, "Ultra-Summarization: A Statistical Approach to Generating Highly Condensed Non-Extractive Summaries", 1999
- [30] Yahoo! <http://www.Yahoo.com>
- [31] Y. Yang and O.J. Pedersen, "A comparative Study o Feature Selection in Text Categorization", Proc. of the fifth International Conference on Machine Learning ICML97, 412-420, 1997